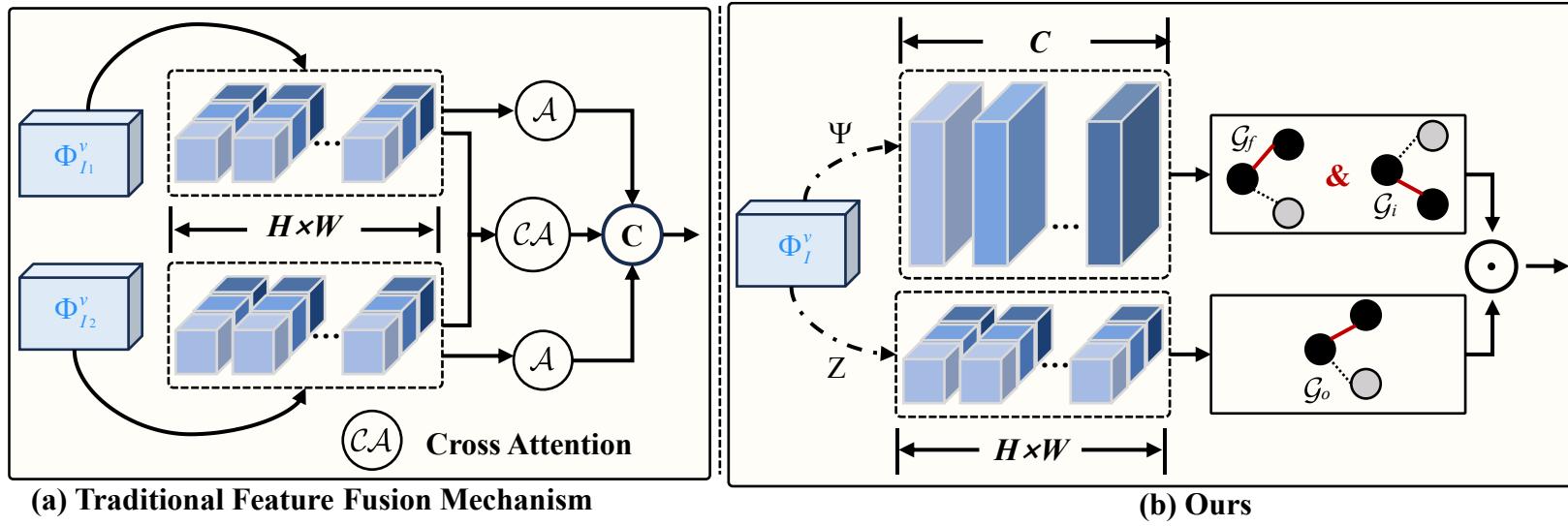
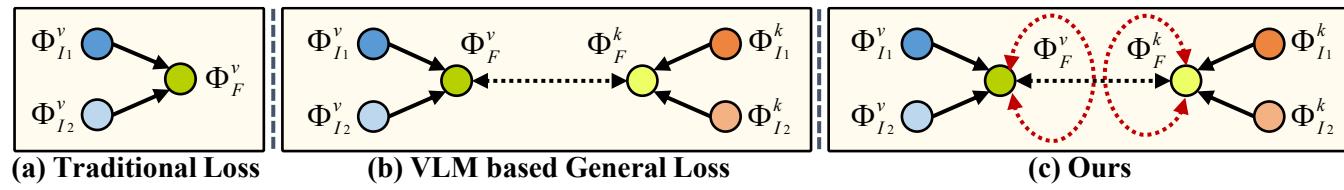
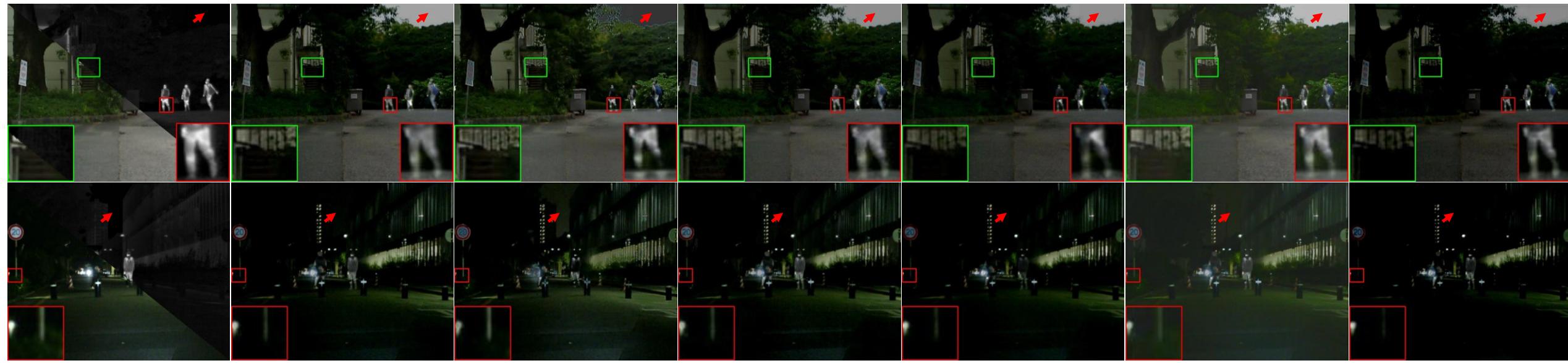


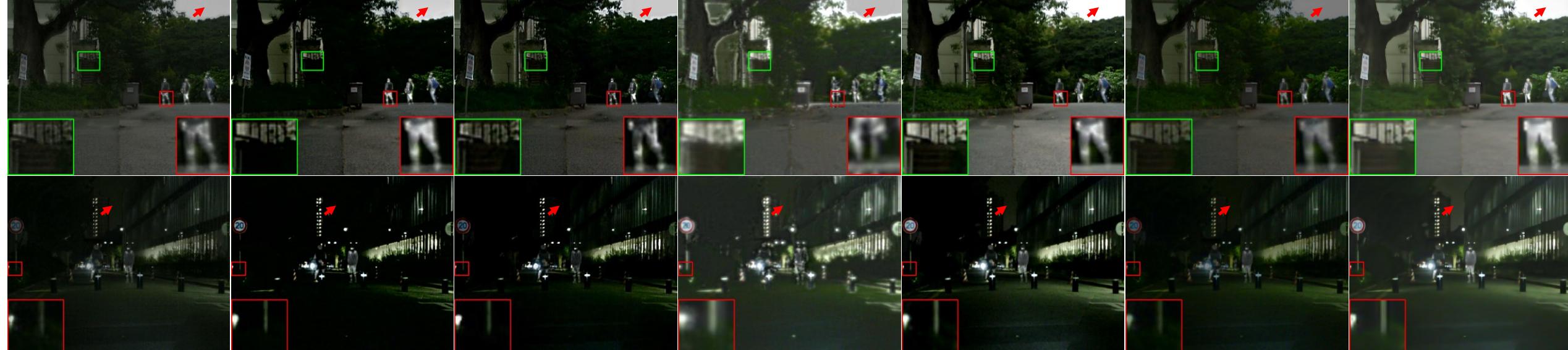
Frozen
 Learnable
 Upsampling
 Down-projection
 Concatenation
 Attention action
 Cell



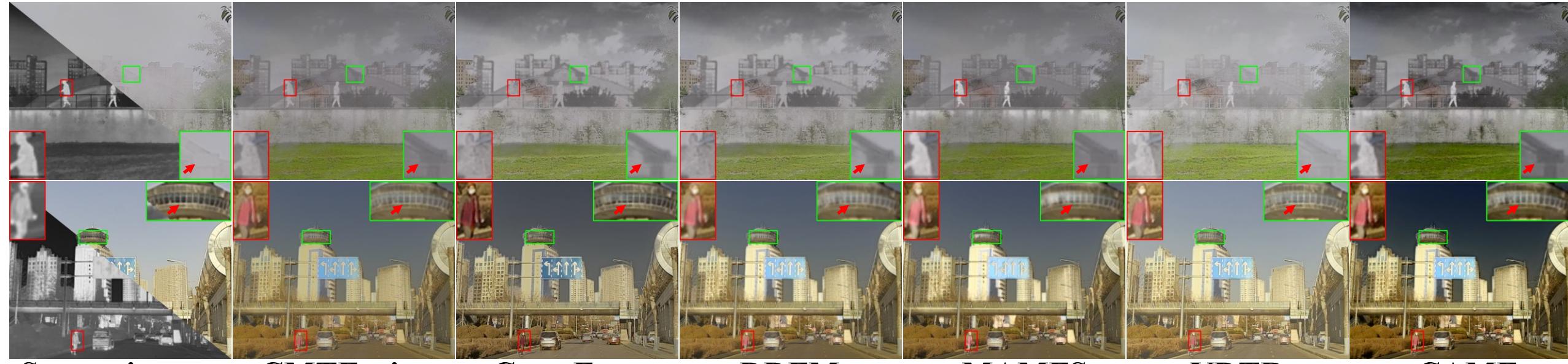




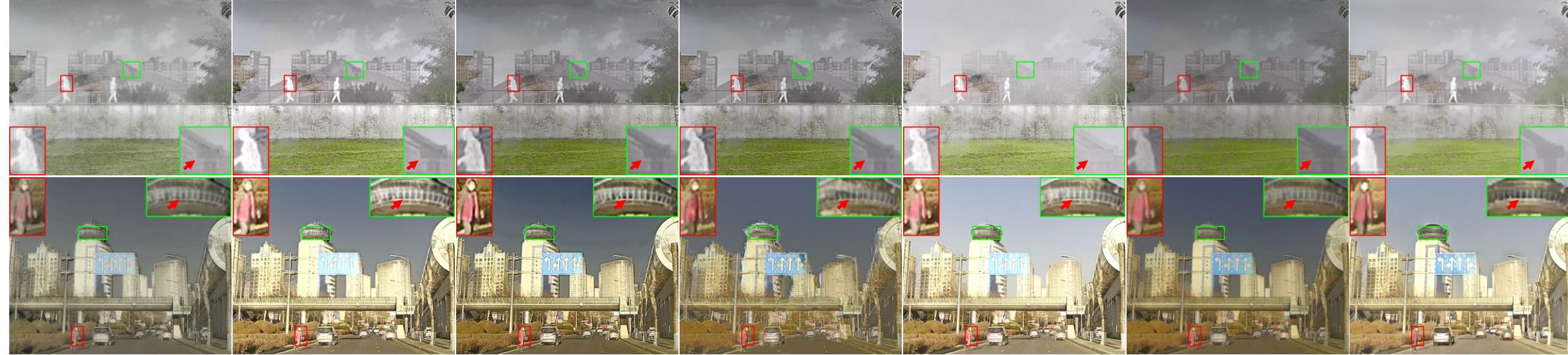
Source images CMTFusion CrossFuse DDFM MAMFS YDTR CAMF



DDBFusion FreeFusion GIFNet MulFS-CAP PromptFusion VDMUFusion Ours



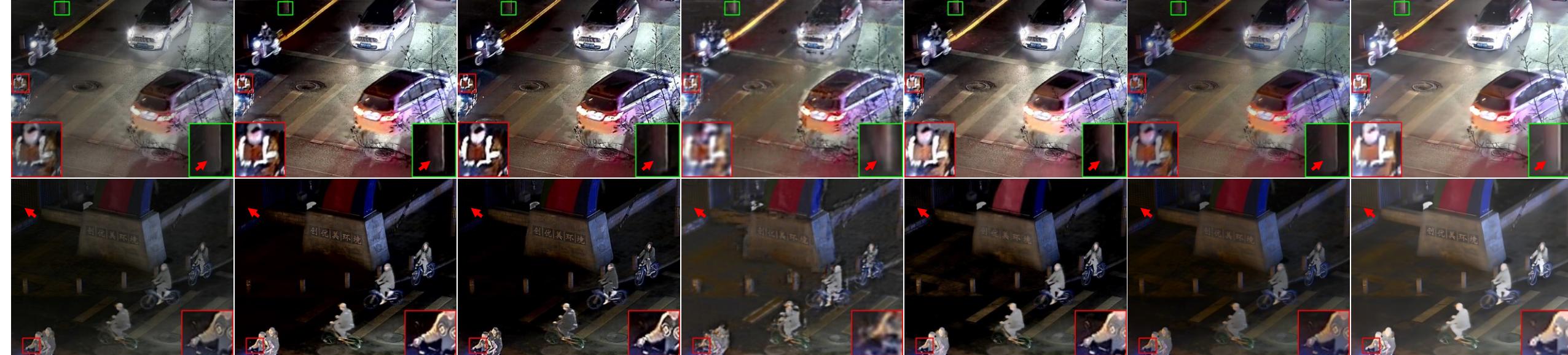
Source images **CMTFusion** **CrossFuse** **DDFM** **MAMFS** **YDTR** **CAMF**



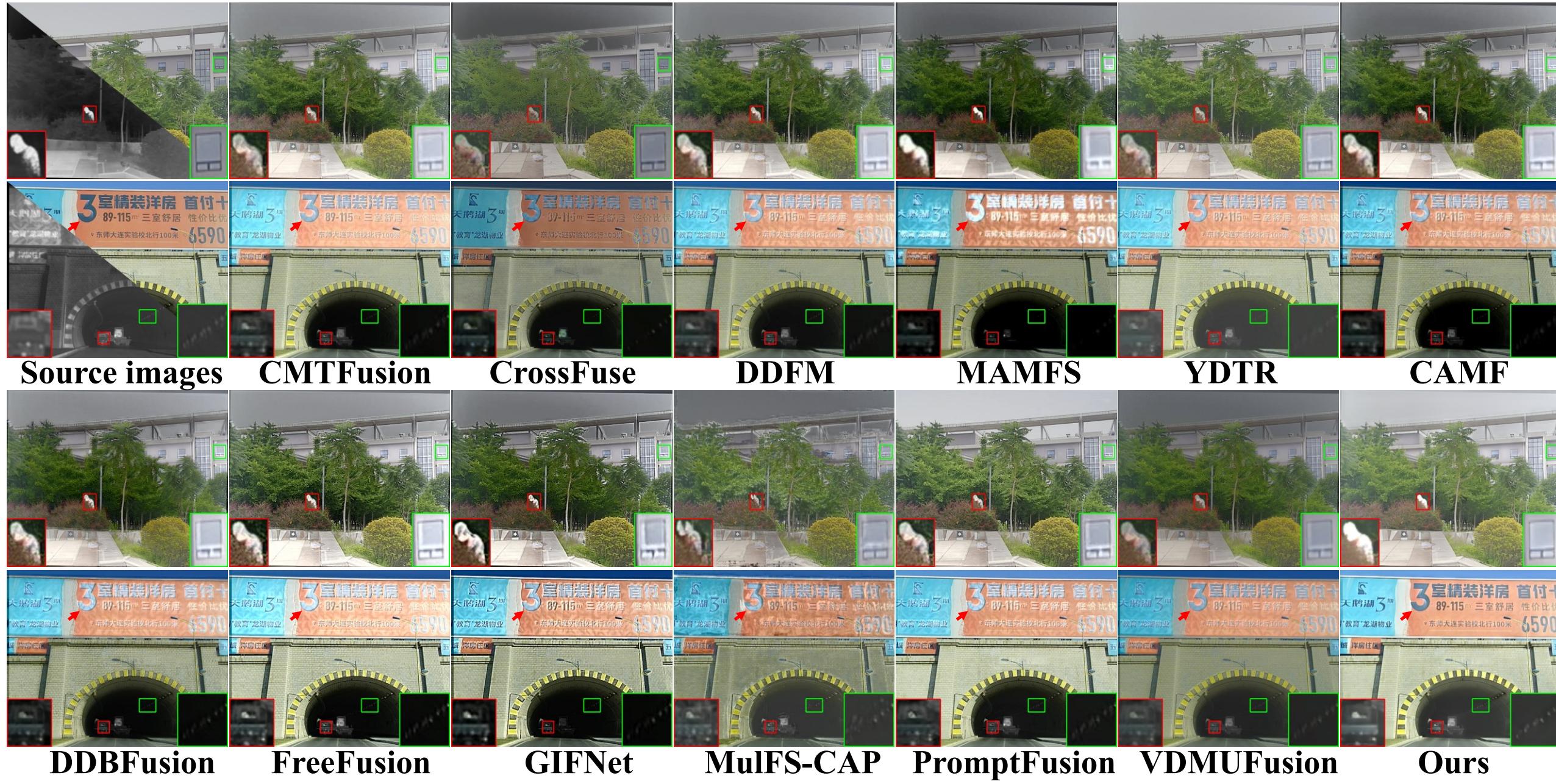
DDBFusion **FreeFusion** **GIFNet** **MulFS-CAP** **PromptFusion** **VDMUFusion** **Ours**

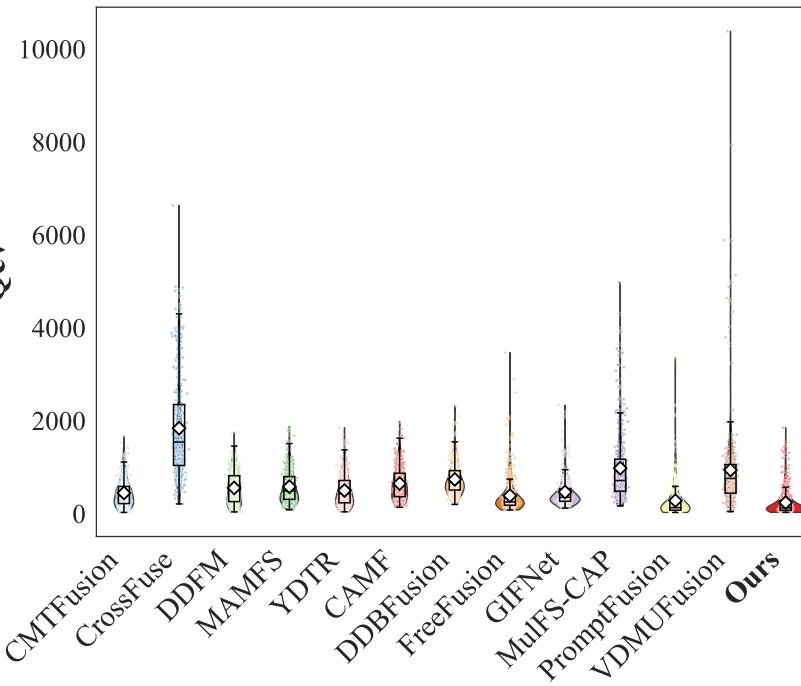
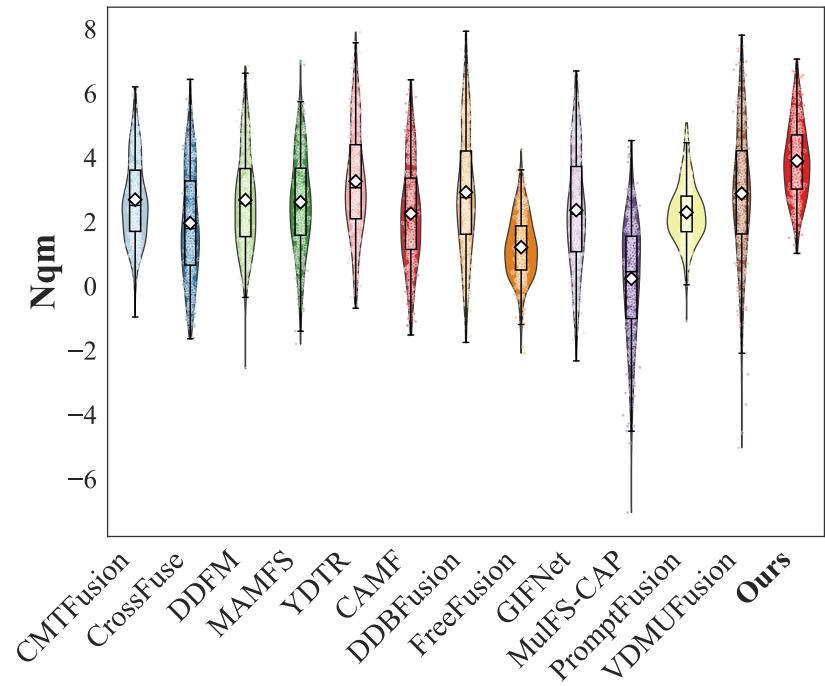
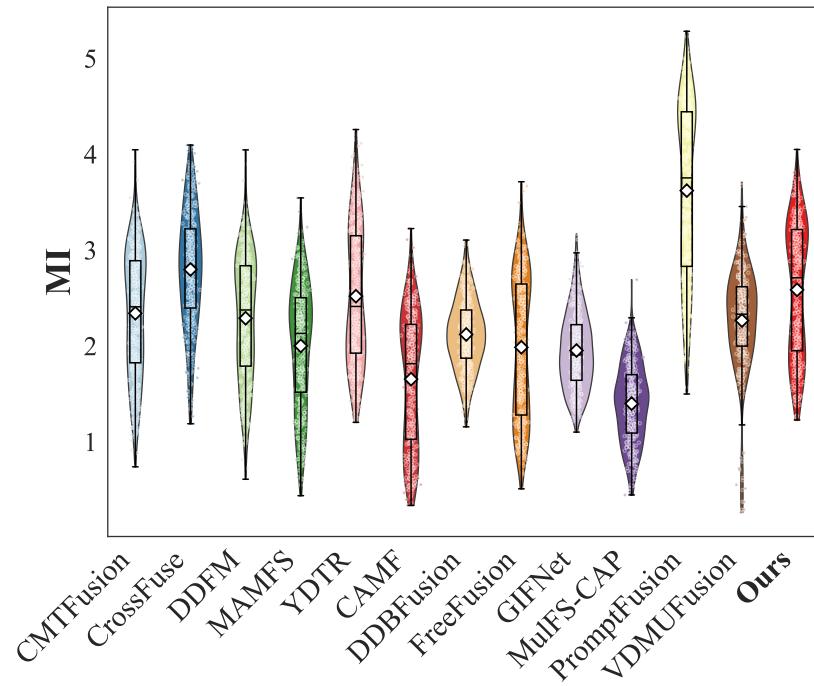
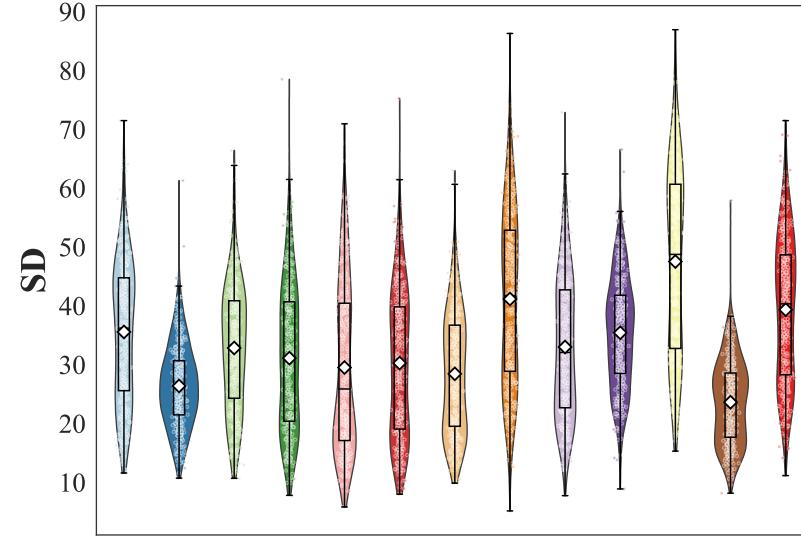
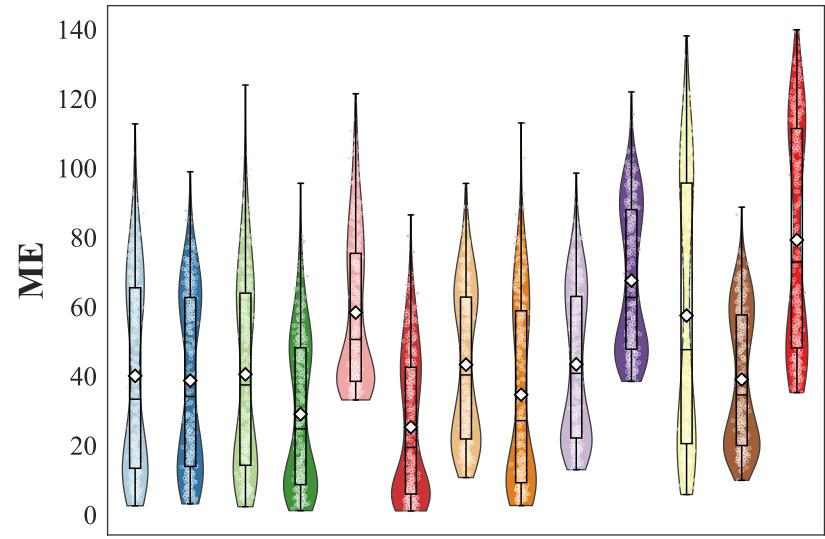


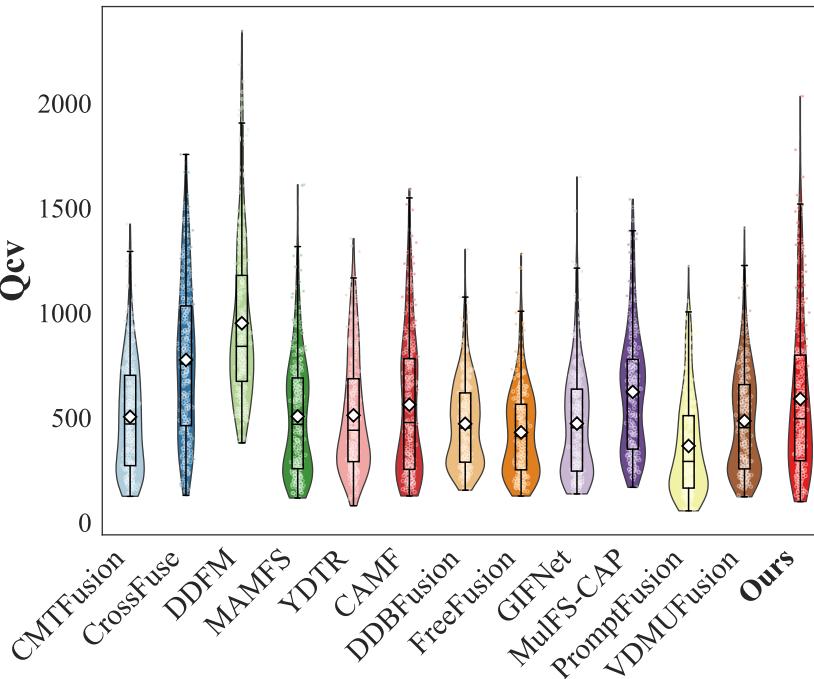
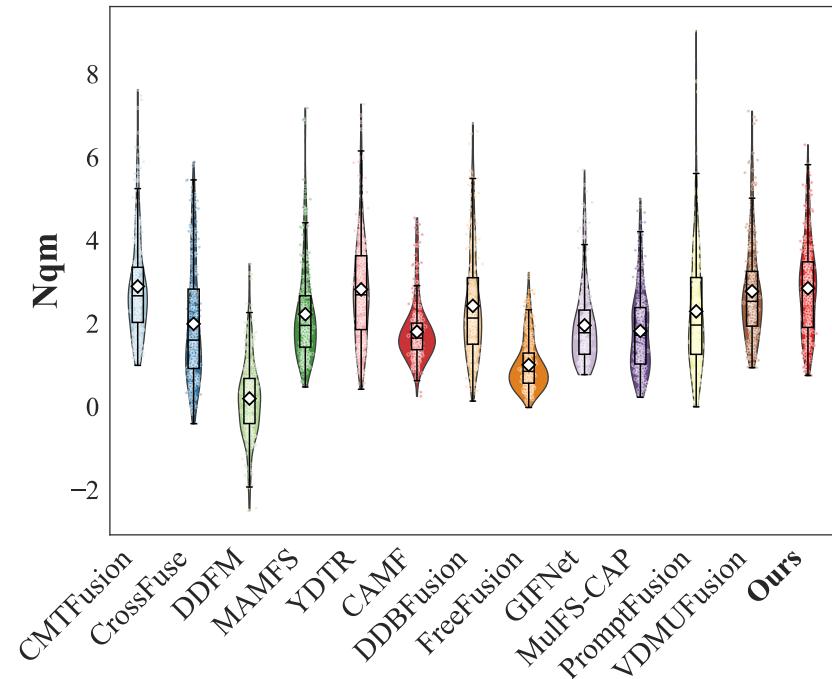
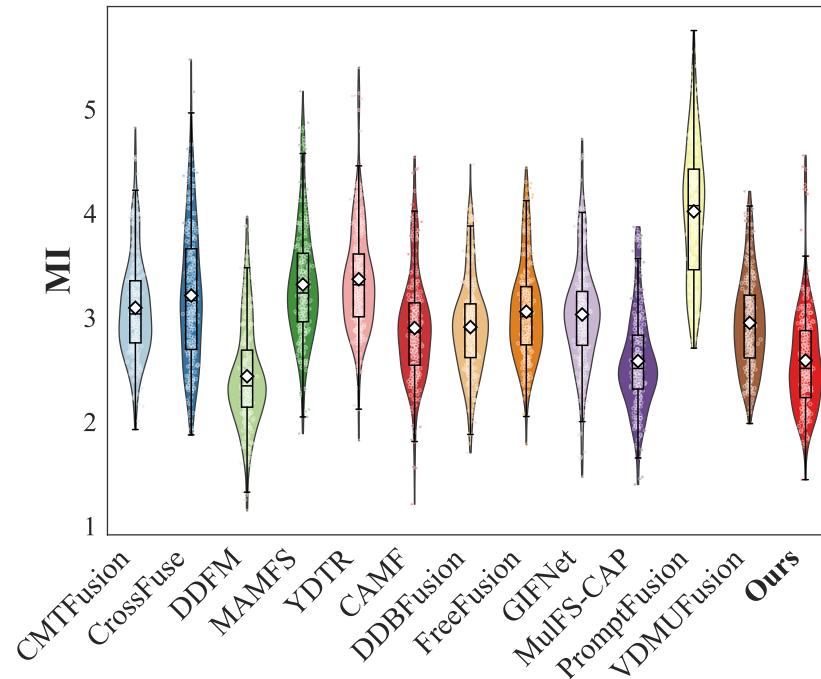
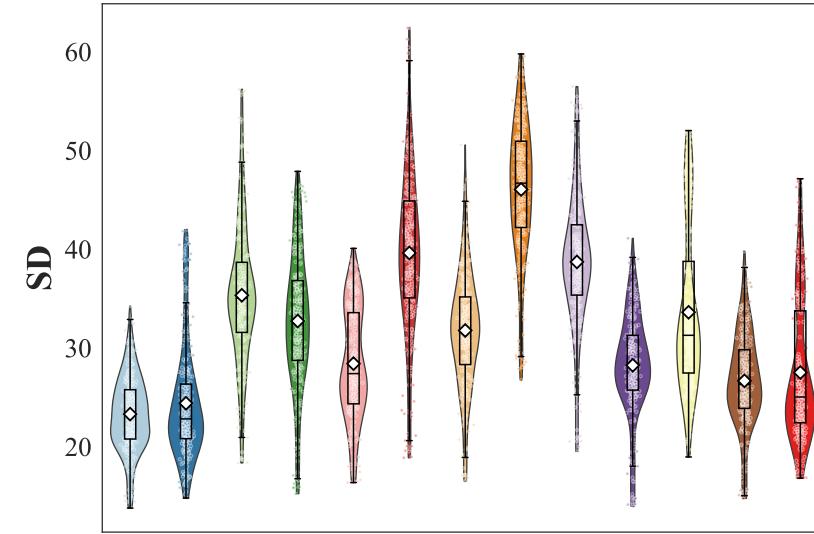
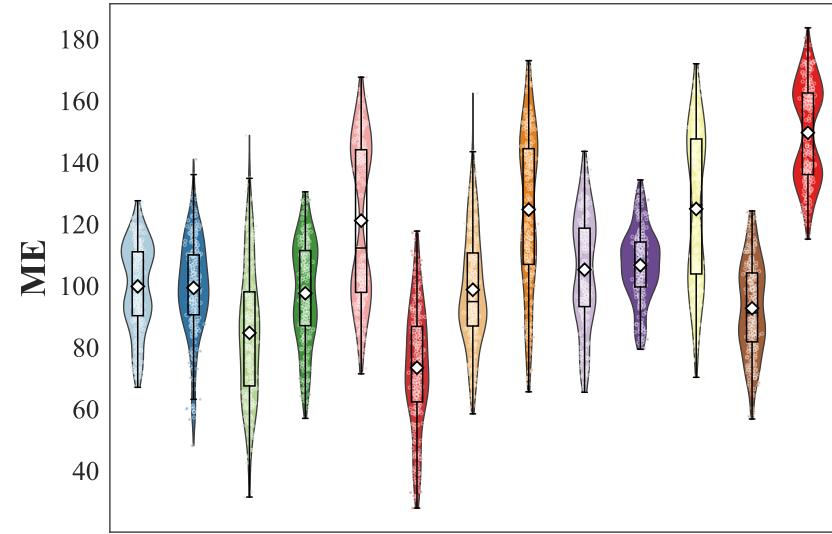
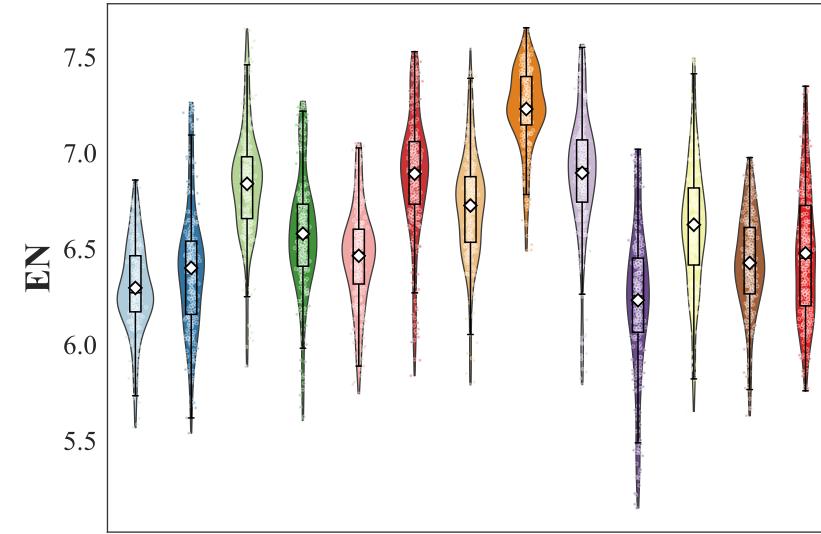
Source images CMTFusion CrossFuse DDFM MAMFS YDTR CAMF

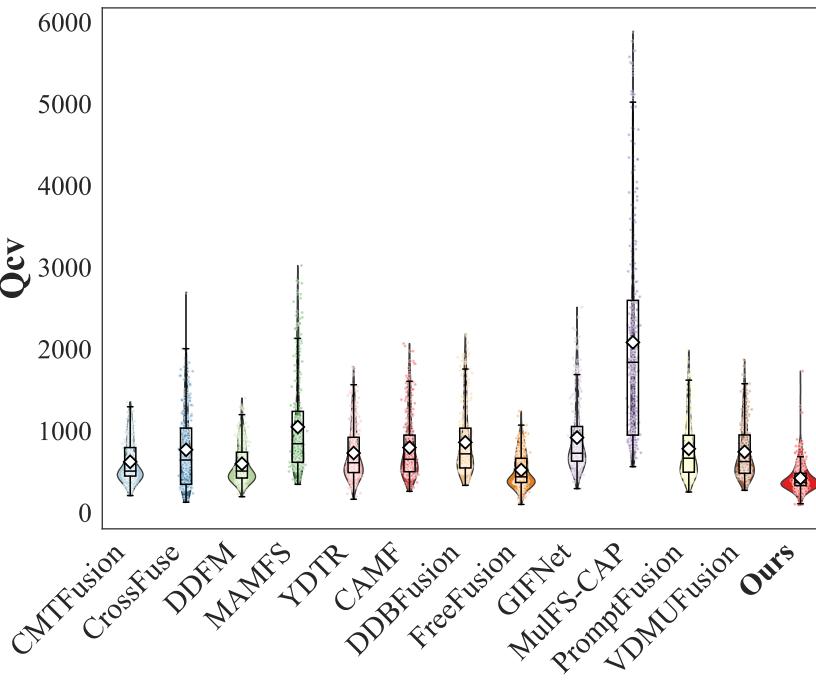
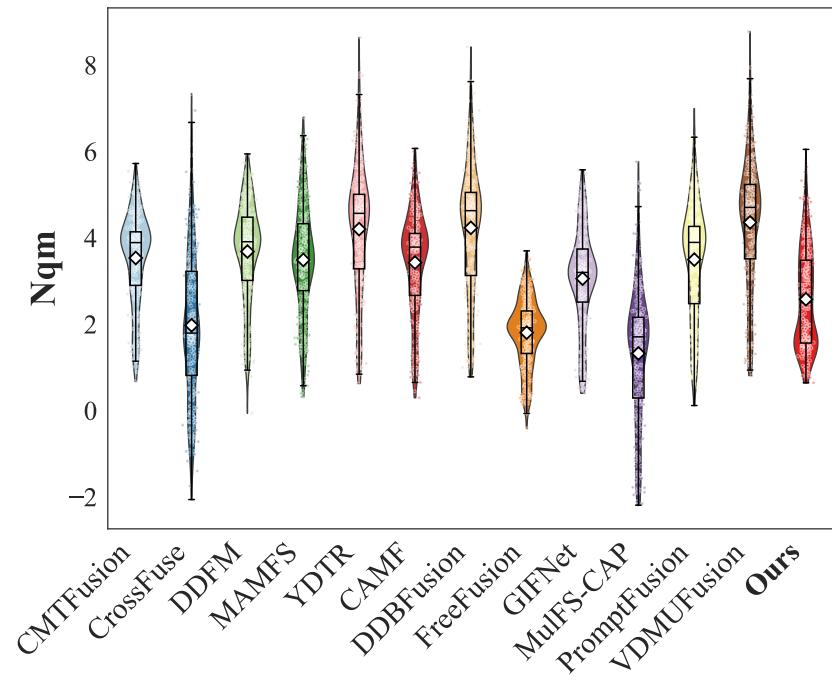
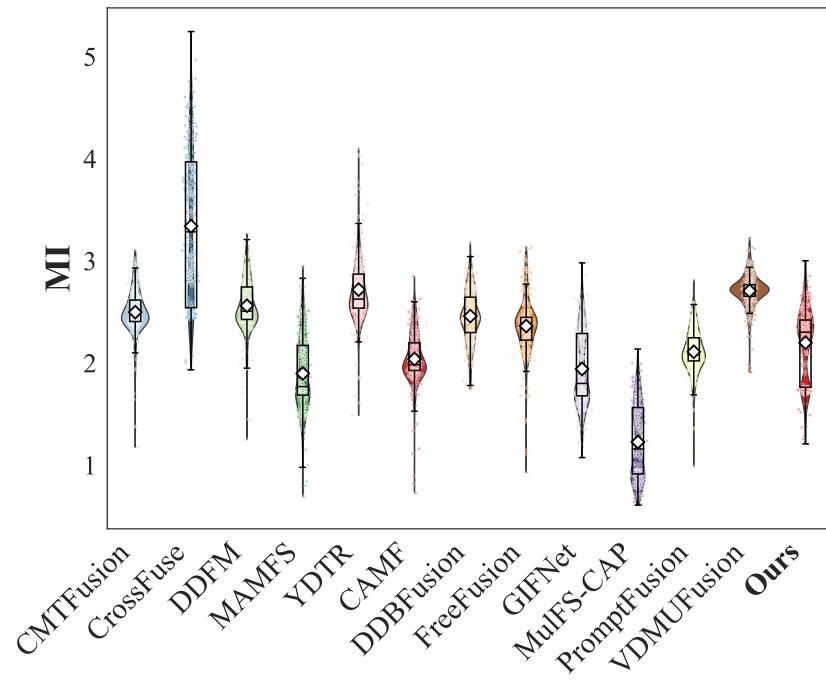
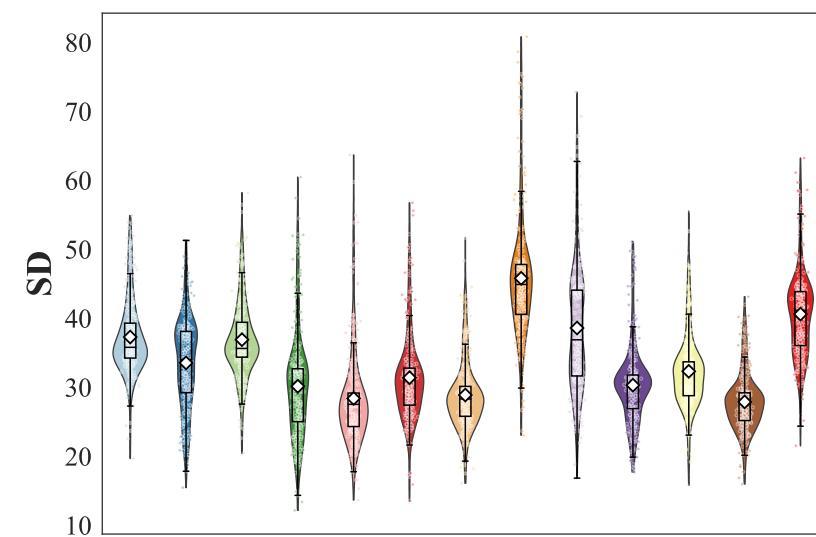
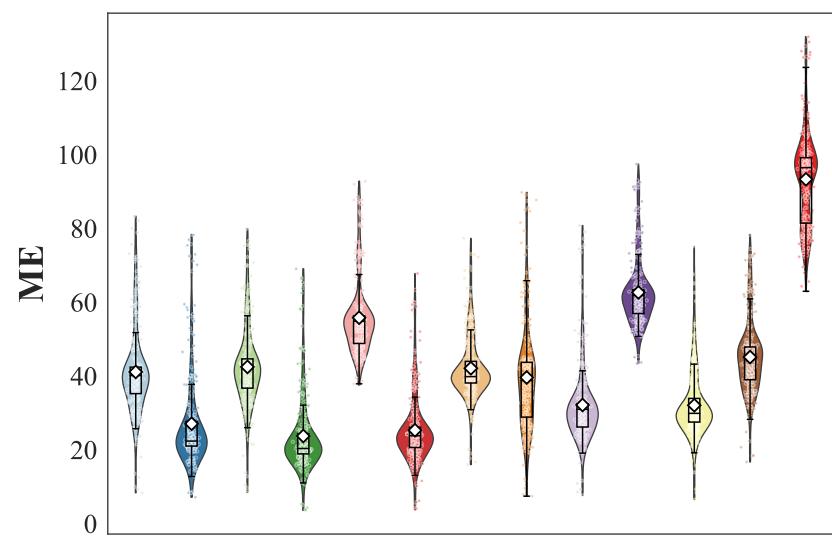
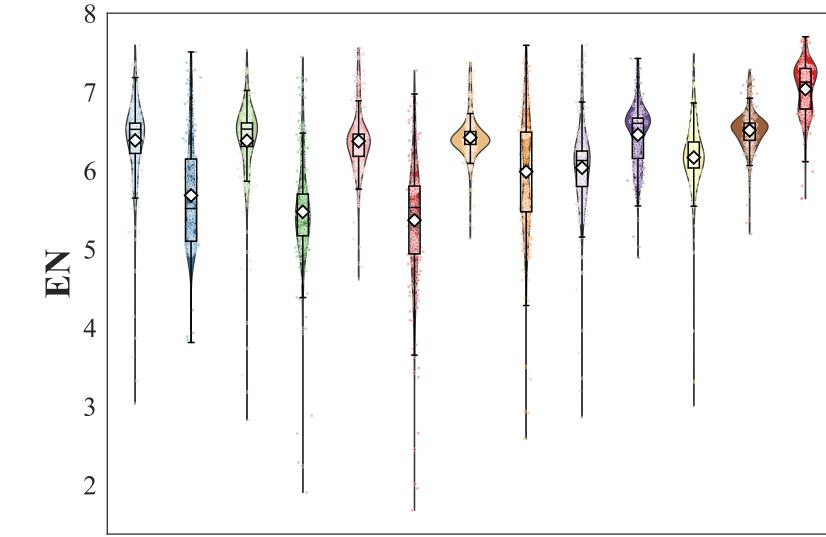


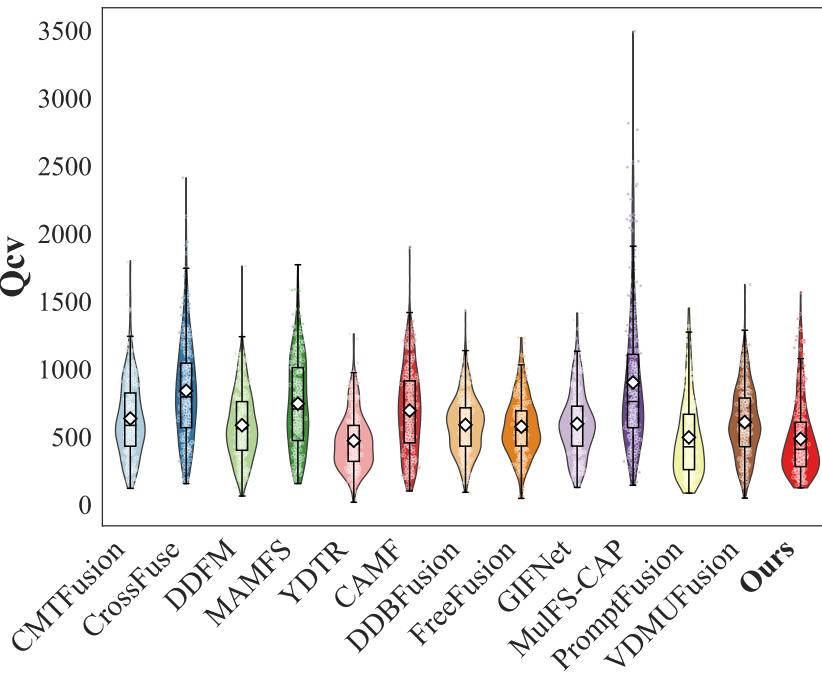
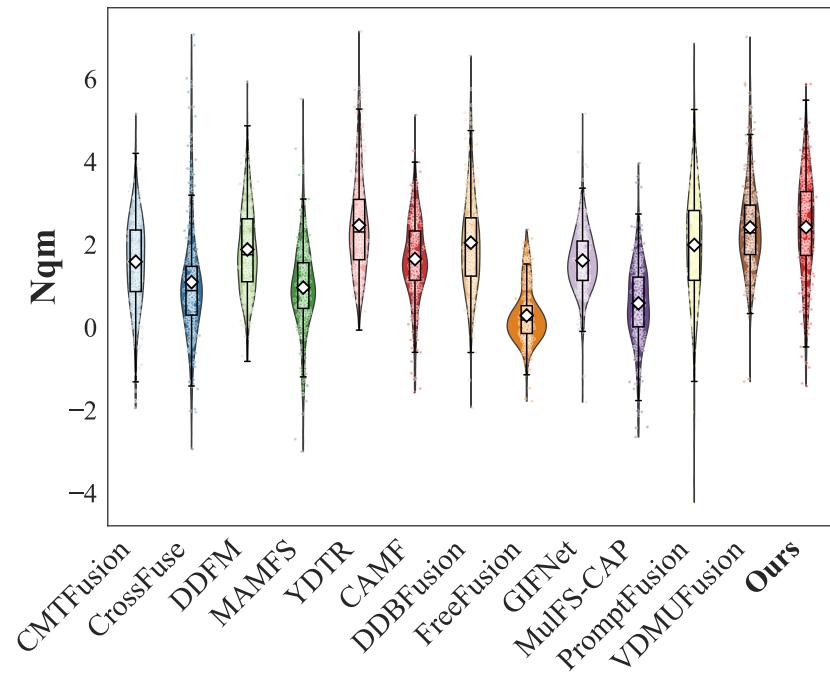
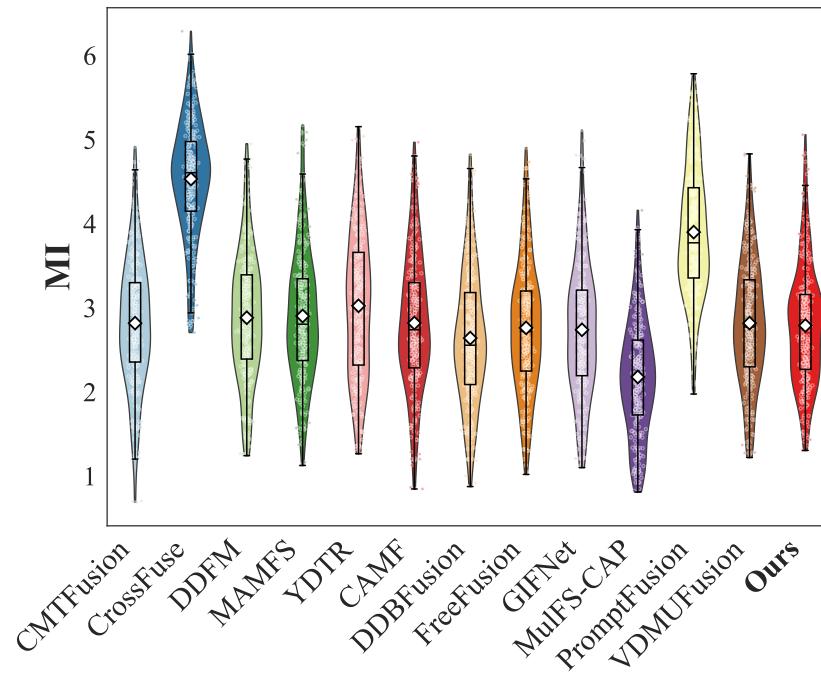
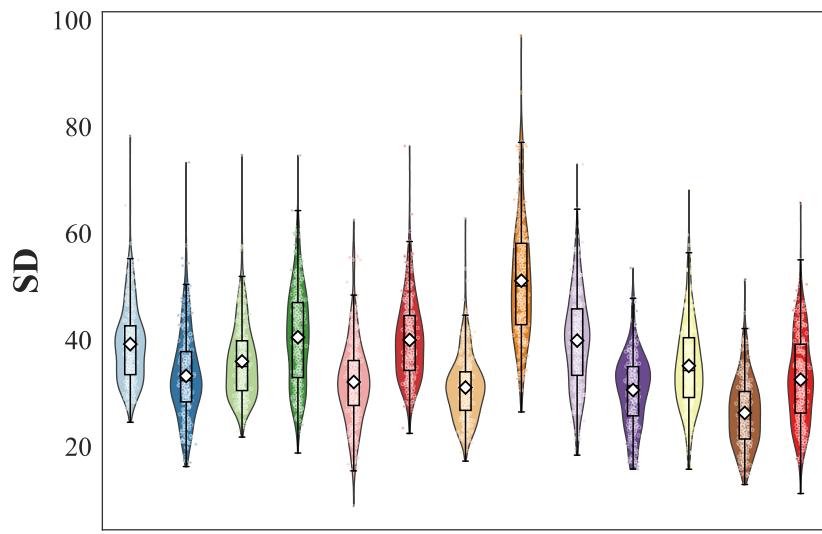
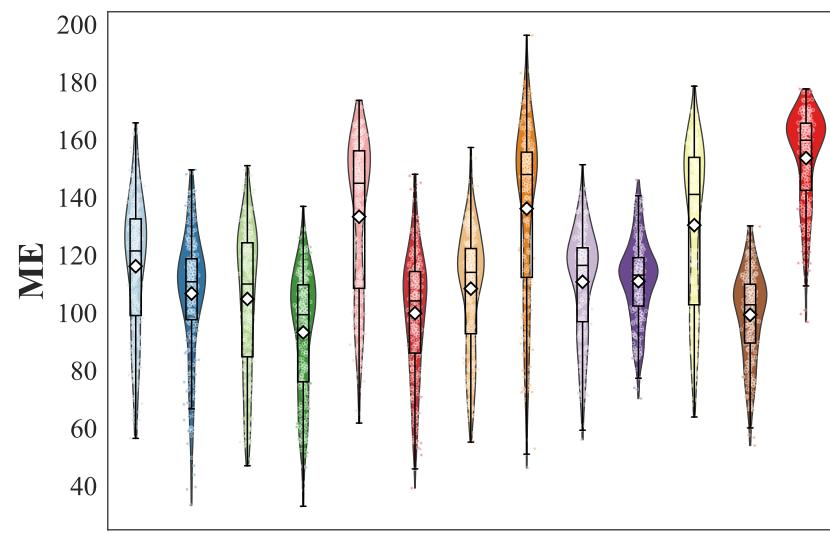
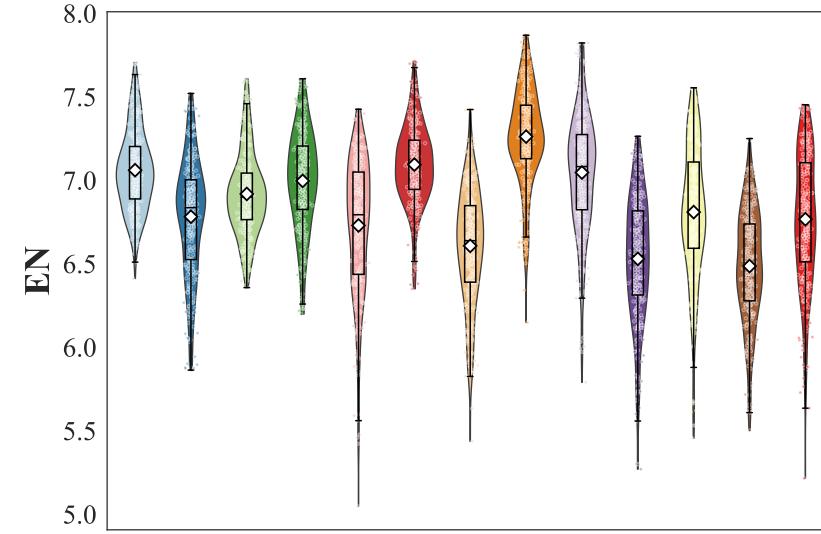
DDBFusion FreeFusion GIFNet MulFS-CAP PromptFusion VDMUFusion Ours

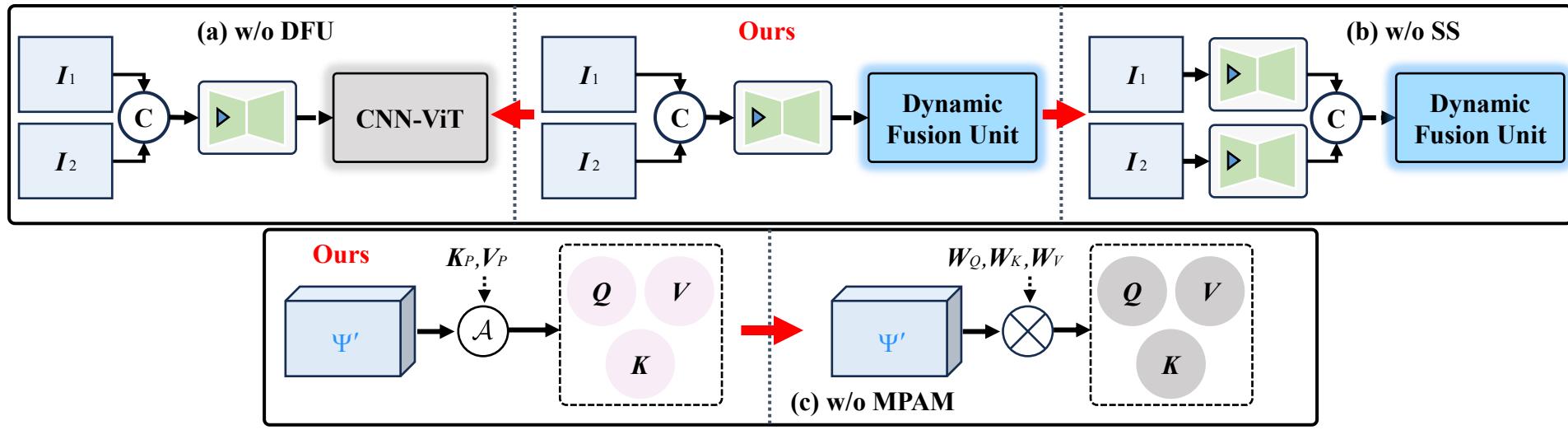


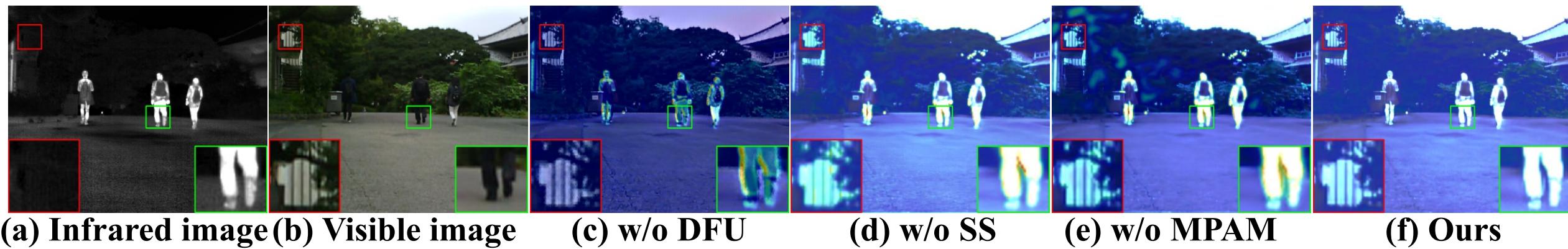














(a) Infrared image

(b) Visible image

(c) w/o \mathcal{L}_s

(d) w/o \mathcal{L}_p



(e) w/o \mathcal{L}_k

(f) $\mathcal{L}_k \rightarrow \mathcal{L}_{VLM_1}$

(g) $\mathcal{L}_k \rightarrow \mathcal{L}_{VLM_2}$

(h) Ours



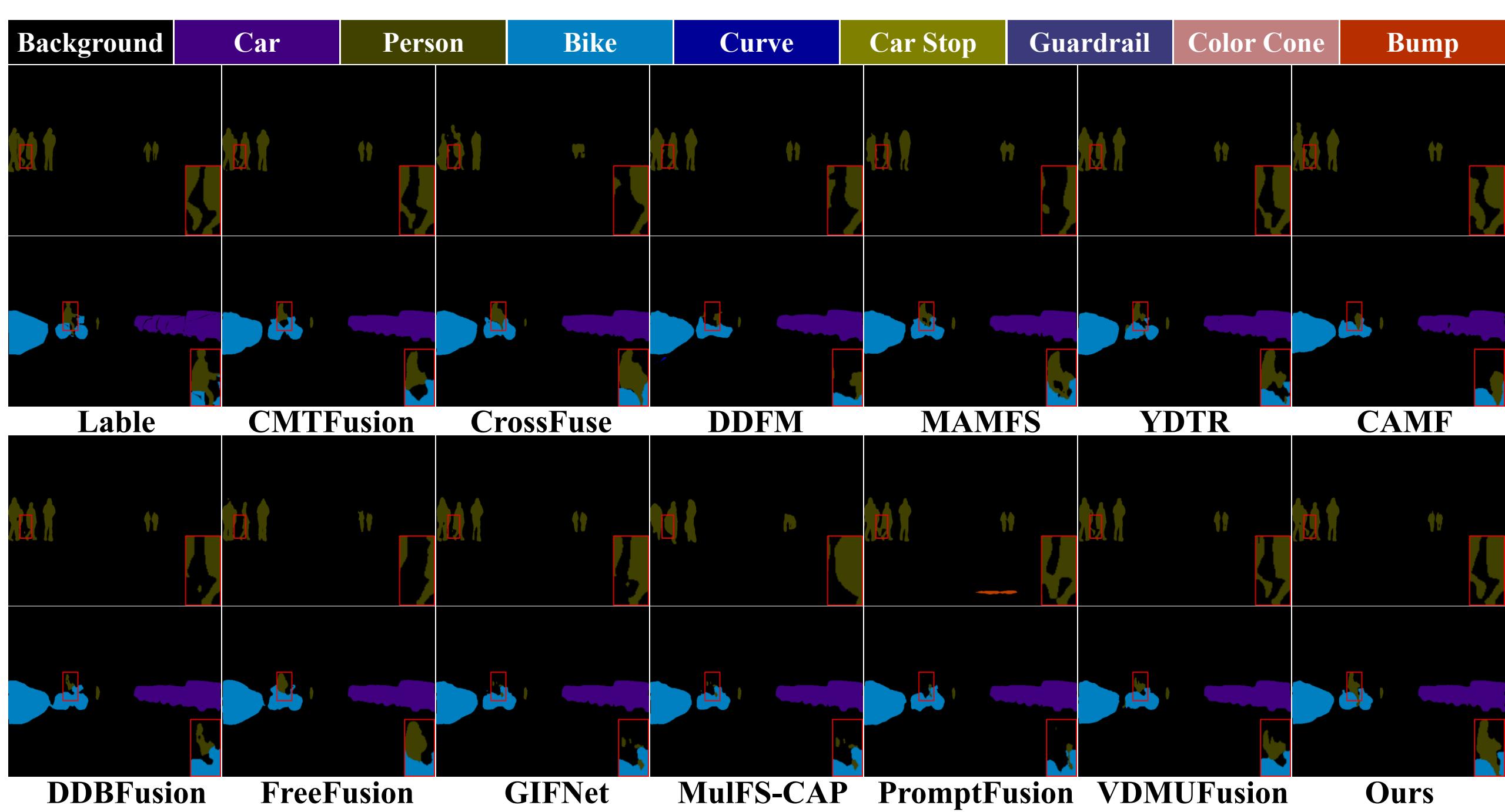
(a) Infrared image (b) Visible image

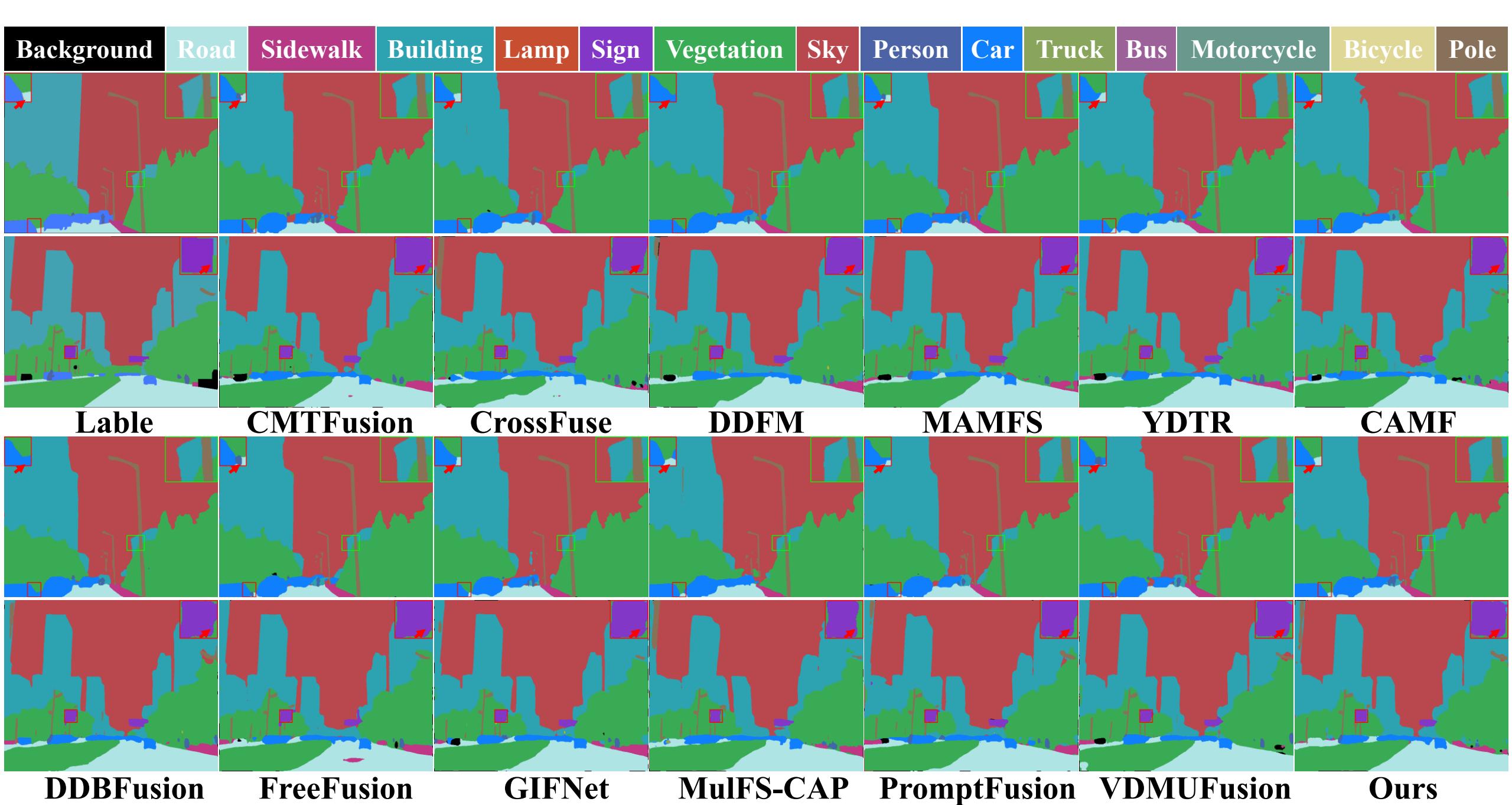
(c) NP

(d) CSP

(e) GFP

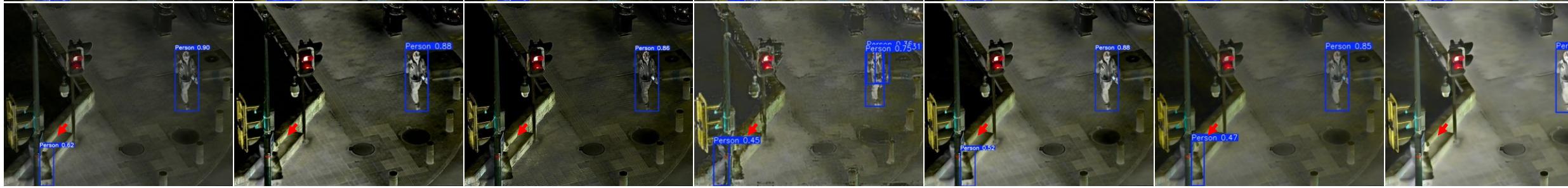
(f) Ours



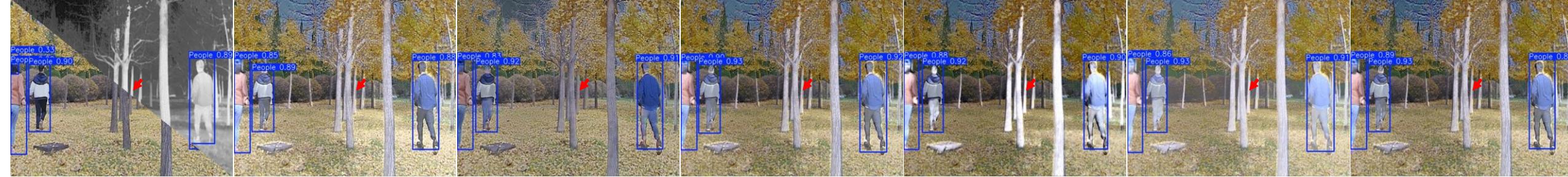




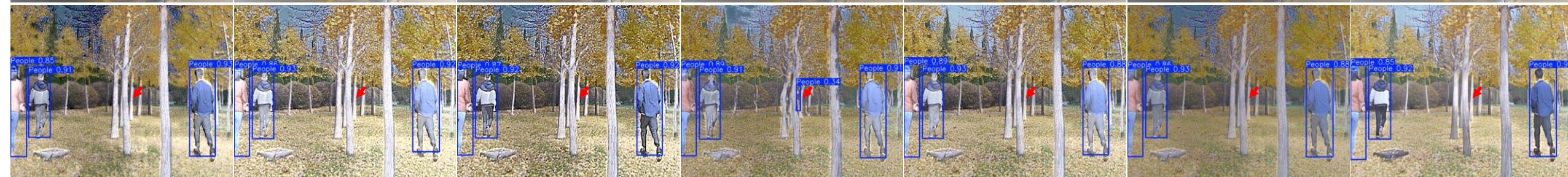
Source images **CMTFusion** **CrossFuse** **DDFM** **MAMFS** **YDTR** **CAMF**



DDBFusion **FreeFusion** **GIFNet** **MulFS-CAP** **PromptFusion** **VDMUFusion** **Ours**



Source images **CMTFusion** **CrossFuse** **DDFM** **MAMFS** **YDTR** **CAMF**



DDBFusion **FreeFusion** **GIFNet** **MulFS-CAP** **PromptFusion** **VDMUFusion** **Ours**