

# On Train-Test Class Overlap and Detection for Image Retrieval

Chull Hwan Song<sup>1</sup> Jooyoung Yoon<sup>1</sup> Taebaek Hwang<sup>1</sup> Shunghyun Choi<sup>1</sup> YeongHyeon Gu<sup>2\*</sup> Yannis Avrithis<sup>3</sup>

<sup>1</sup>Dealicious Inc. <sup>2</sup>Sejong University <sup>3</sup>Institute of Advanced Research on Artificial Intelligence (IARAI)



IARAI

institute of advanced  
research in artificial  
intelligence

# Introduction

- **Background:** Importance of non-overlapping training and evaluation sets in image retrieval.
- **Problem Statement:** Existing methods do not adequately address class overlap issues in datasets.
- **Objective:** Introduce RGLDv2-clean and CiDeR for effective image retrieval.

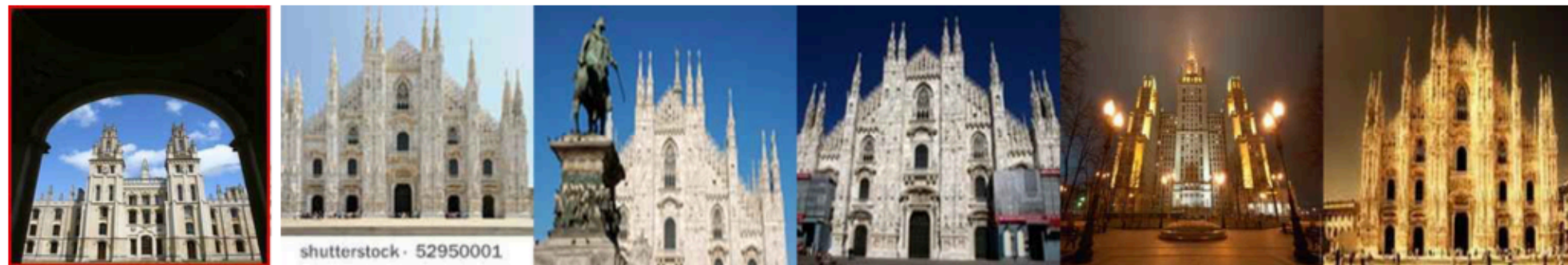


# Main Contributions

Training:  
GLDv2-clean



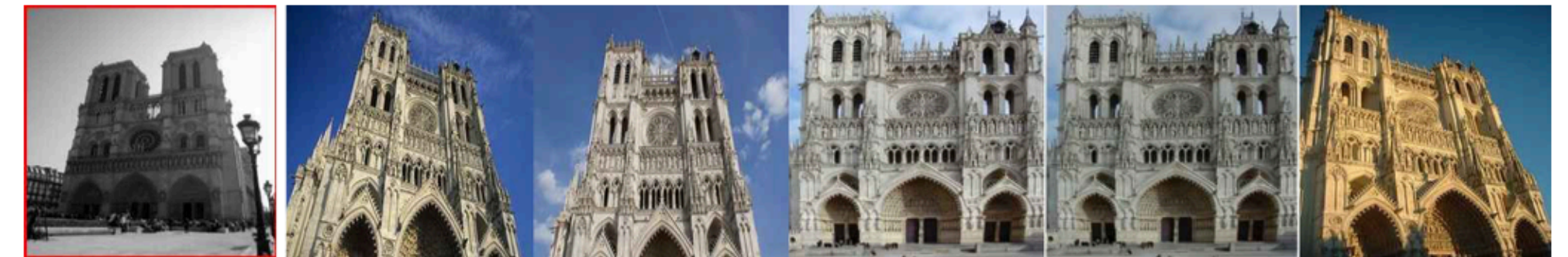
Training:  
NC-clean



Training:  
SfM-120k



Evaluation:  $\mathcal{R}$ Oxford



Evaluation:  $\mathcal{R}$ Paris

- Key Points: Overlap between GLDv2-clean and evaluation sets( $\mathcal{R}$ Oxford and  $\mathcal{R}$ Par)

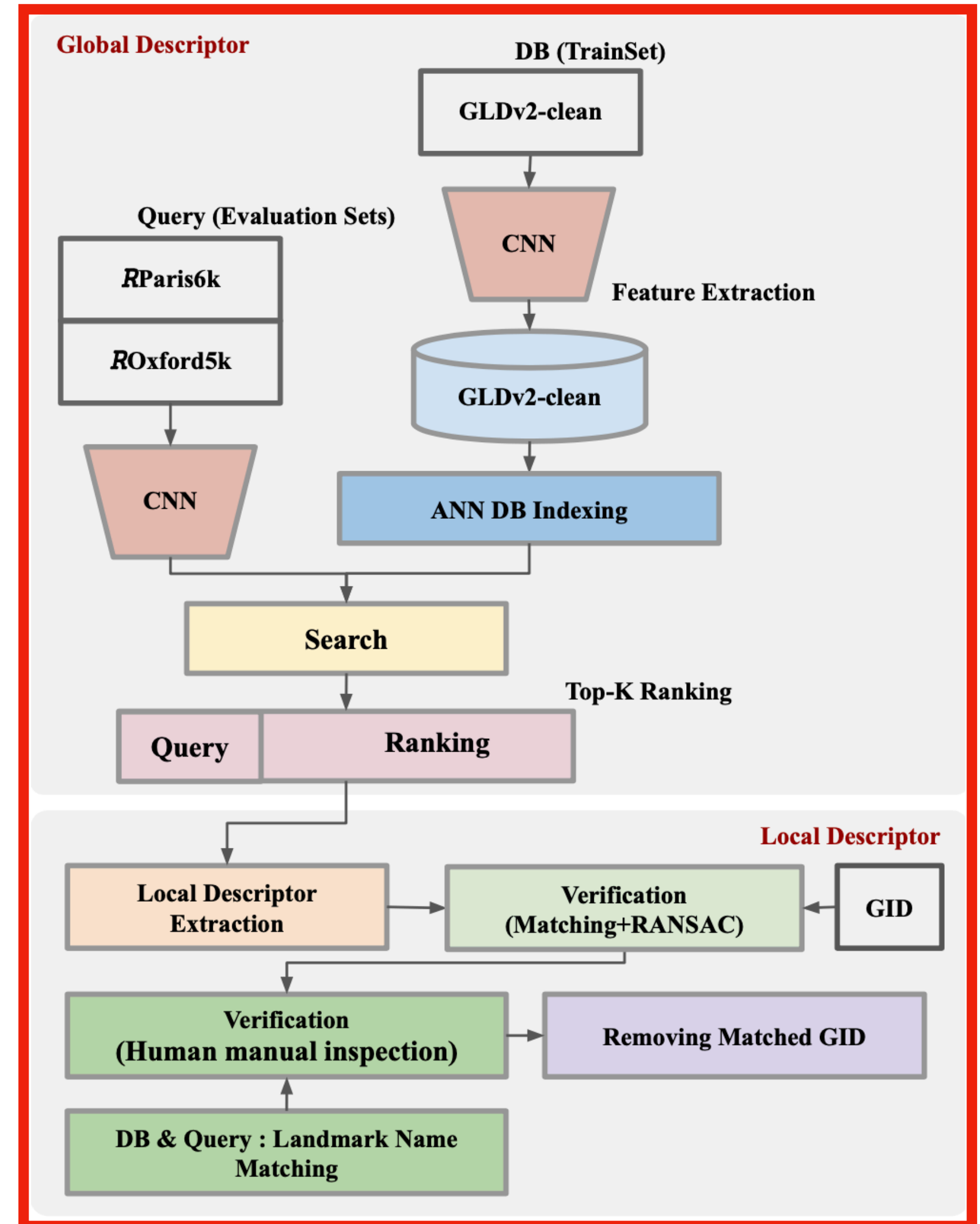


# Data Cleaning Process

- Steps: Identify, remove overlaps, verify
- Table: Statistics before and after cleaning

EVAL	#EVAL IMG	#DUPL EVAL	#DUPL GLDV2	GID	#DUPL GLDV2 IMG
$\mathcal{R}Par$	70	36 (51%)	11		1,227
$\mathcal{R}Oxf$	70	38 (54%)	6		315
TEXT			1		23
TOTAL	140	74	18		1,565

TRAINING SET	#IMAGES	#CATEGORIES
NC-clean	27,965	581
SfM-120k	117,369	713
GLDv2-clean	1,580,470	81,313
$\mathcal{R}GLDv2$ -clean (ours)	1,578,905	81,295

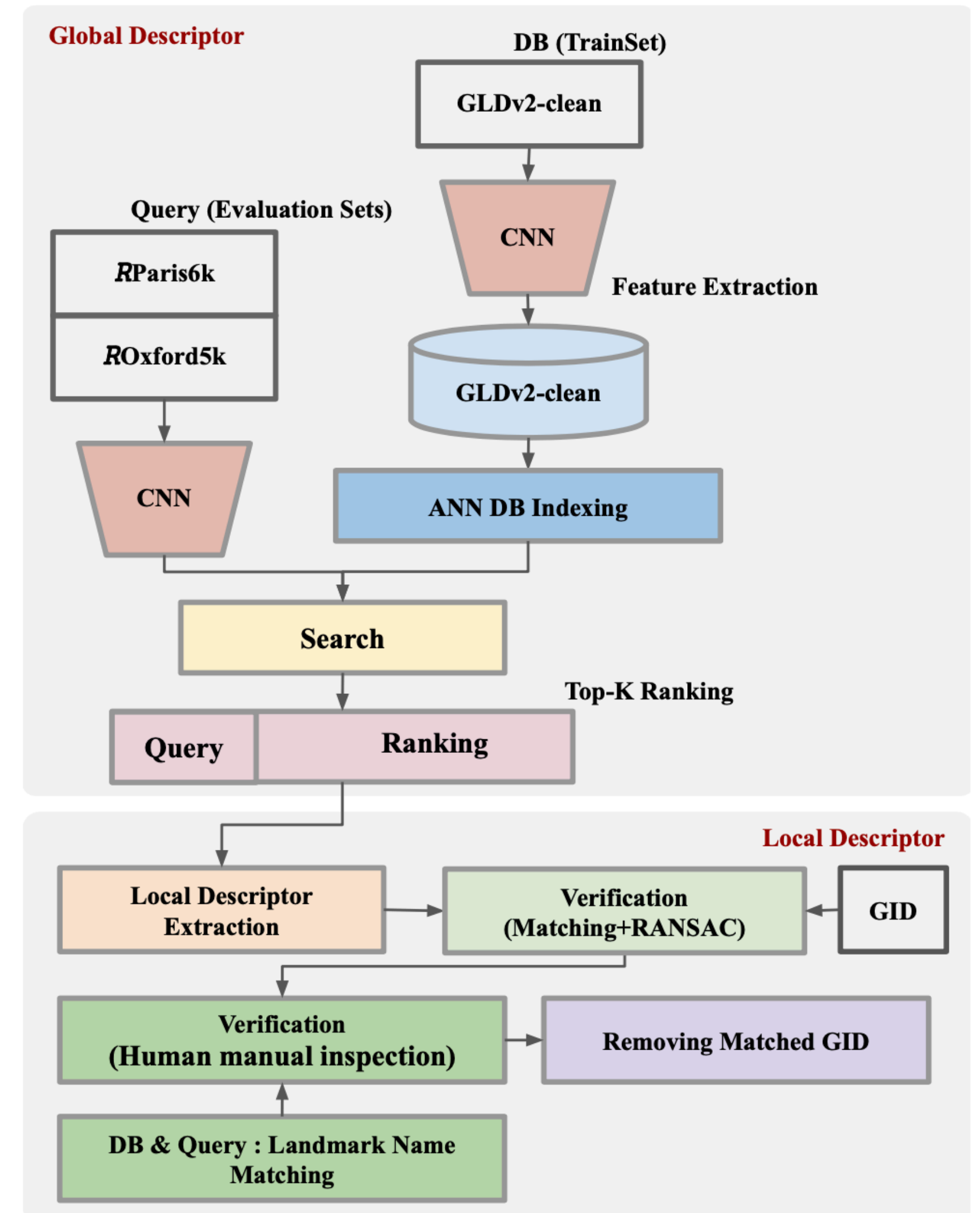


# Data Cleaning Process

- Steps: Identify, remove overlaps, verify
- Table: Statistics before and after cleaning

EVAL	#EVAL IMG	#DUPL EVAL	#DUPL GLDV2	GID	#DUPL GLDV2 IMG
$\mathcal{R}Par$	70	36 (51%)	11		1,227
$\mathcal{R}Oxf$	70	38 (54%)	6		315
TEXT			1		23
TOTAL	140	74	18		1,565

TRAINING SET	#IMAGES	#CATEGORIES
NC-clean	27,965	581
SfM-120k	117,369	713
GLDv2-clean	1,580,470	81,313
$\mathcal{R}GLDv2$ -clean (ours)	1,578,905	81,295

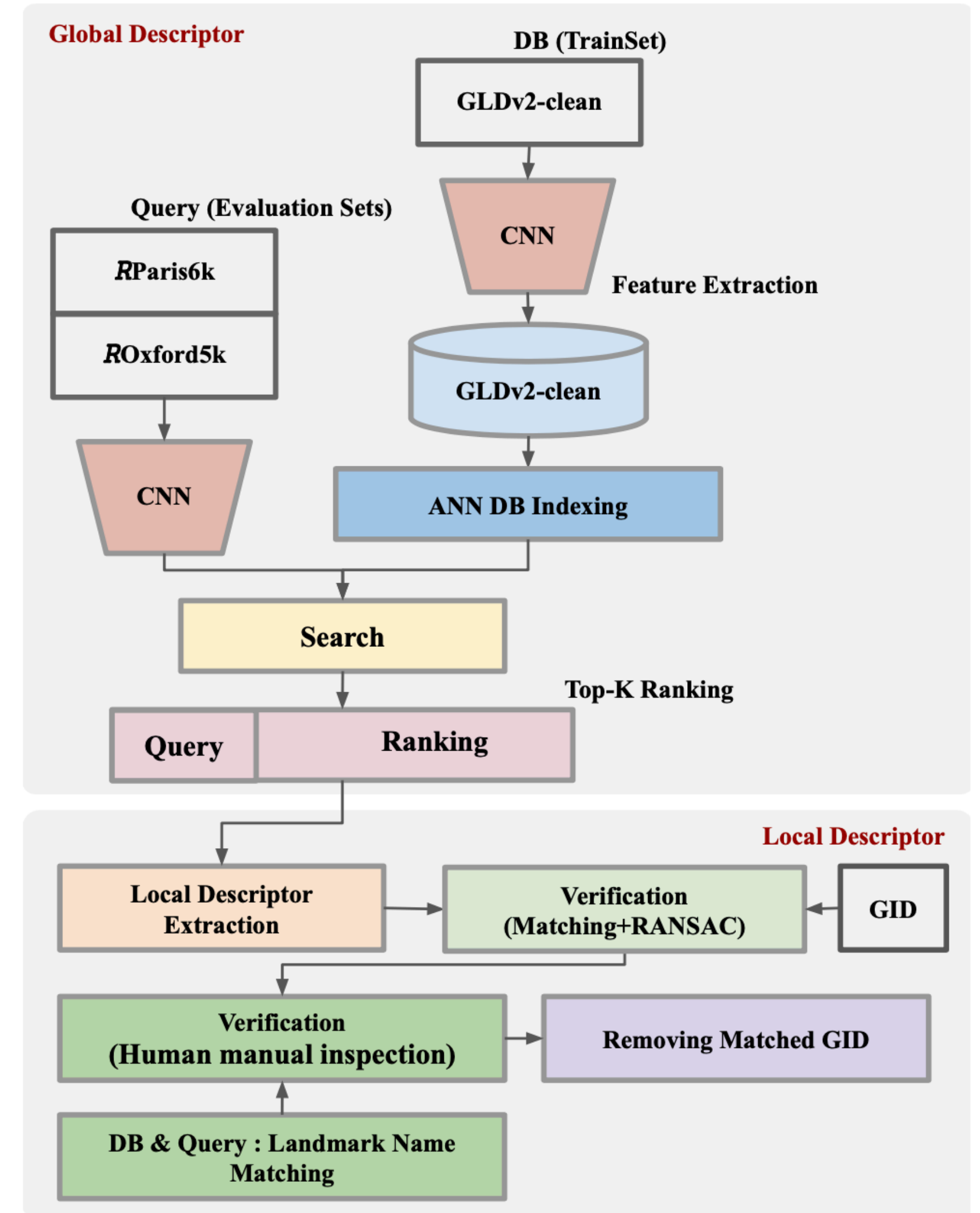


# Data Cleaning Process

- Steps: Identify, remove overlaps, verify
- Table: Statistics before and after cleaning

EVAL	#EVAL IMG	#DUPL EVAL	#DUPL GLDV2	GID	#DUPL GLDV2 IMG
$\mathcal{R}Par$	70	36 (51%)	11		1,227
$\mathcal{R}Oxf$	70	38 (54%)	6		315
TEXT			1		23
TOTAL	140	74	18		1,565

TRAINING SET	#IMAGES	#CATEGORIES
NC-clean	27,965	581
SfM-120k	117,369	713
GLDv2-clean	1,580,470	81,313
$\mathcal{R}GLDv2$ -clean (ours)	1,578,905	81,295



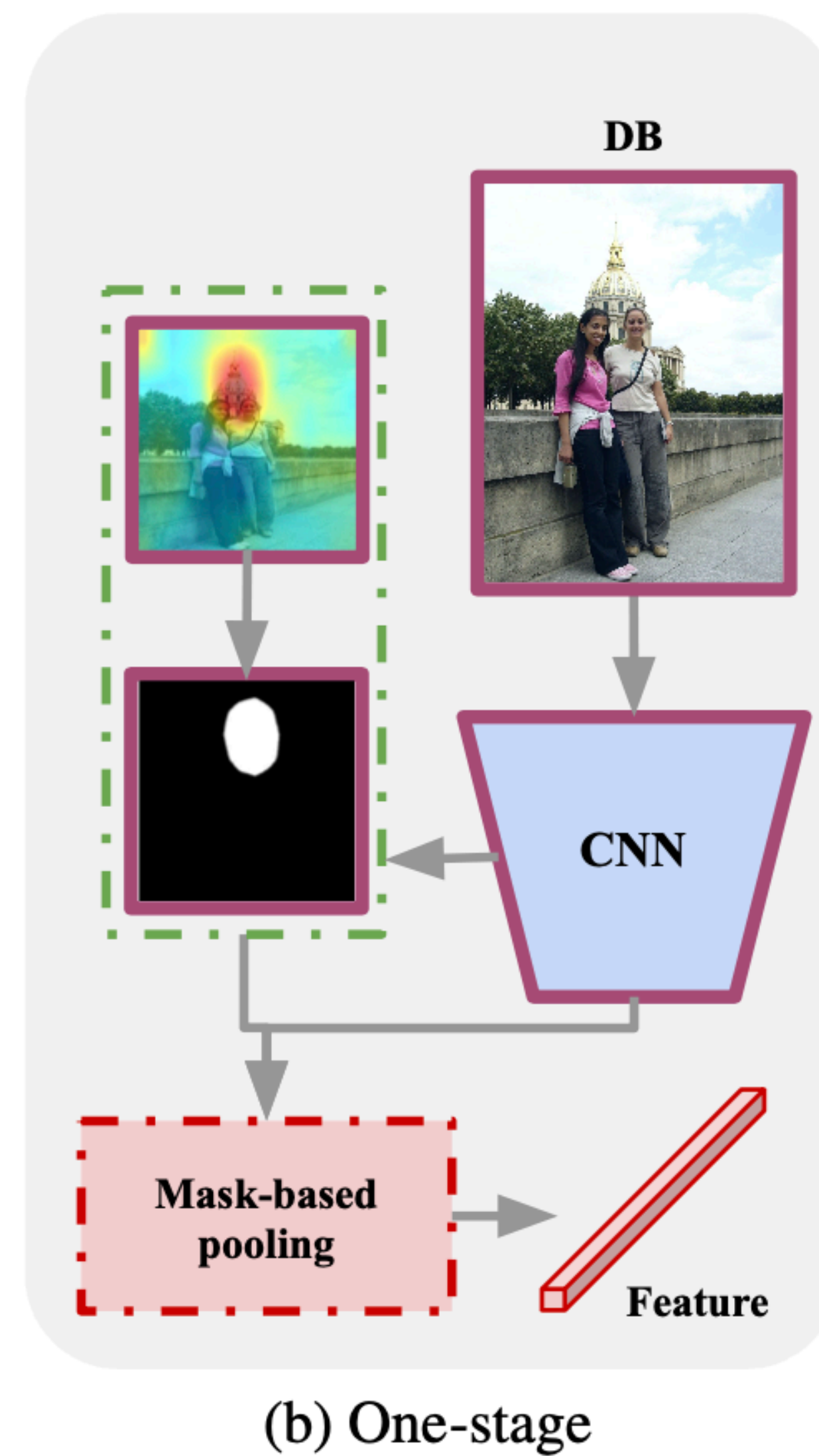
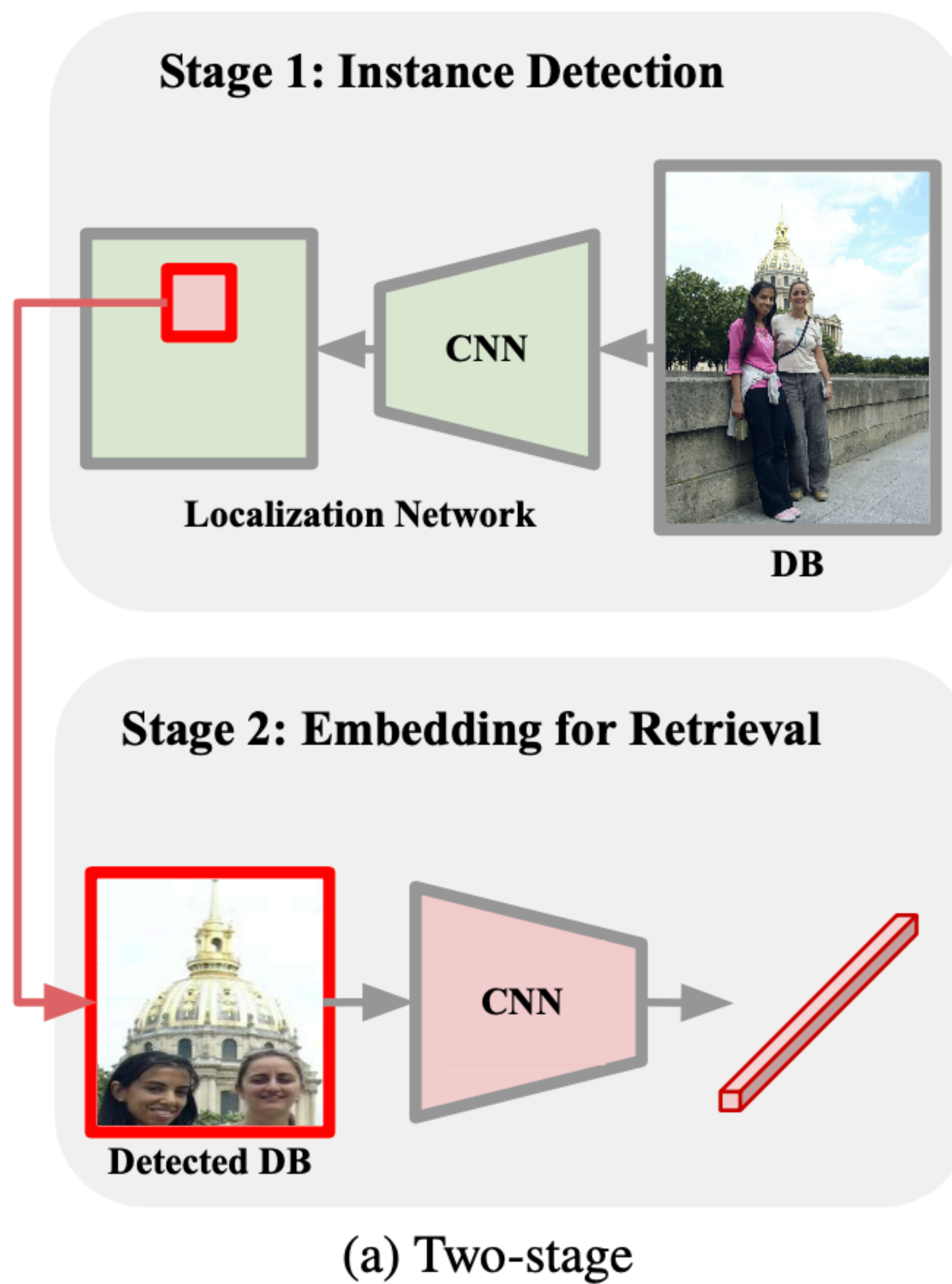


# Performance Impact

- Performance comparison on GLDv2-clean vs RGLDv2-clean
- Significant drop in performance on cleaned dataset

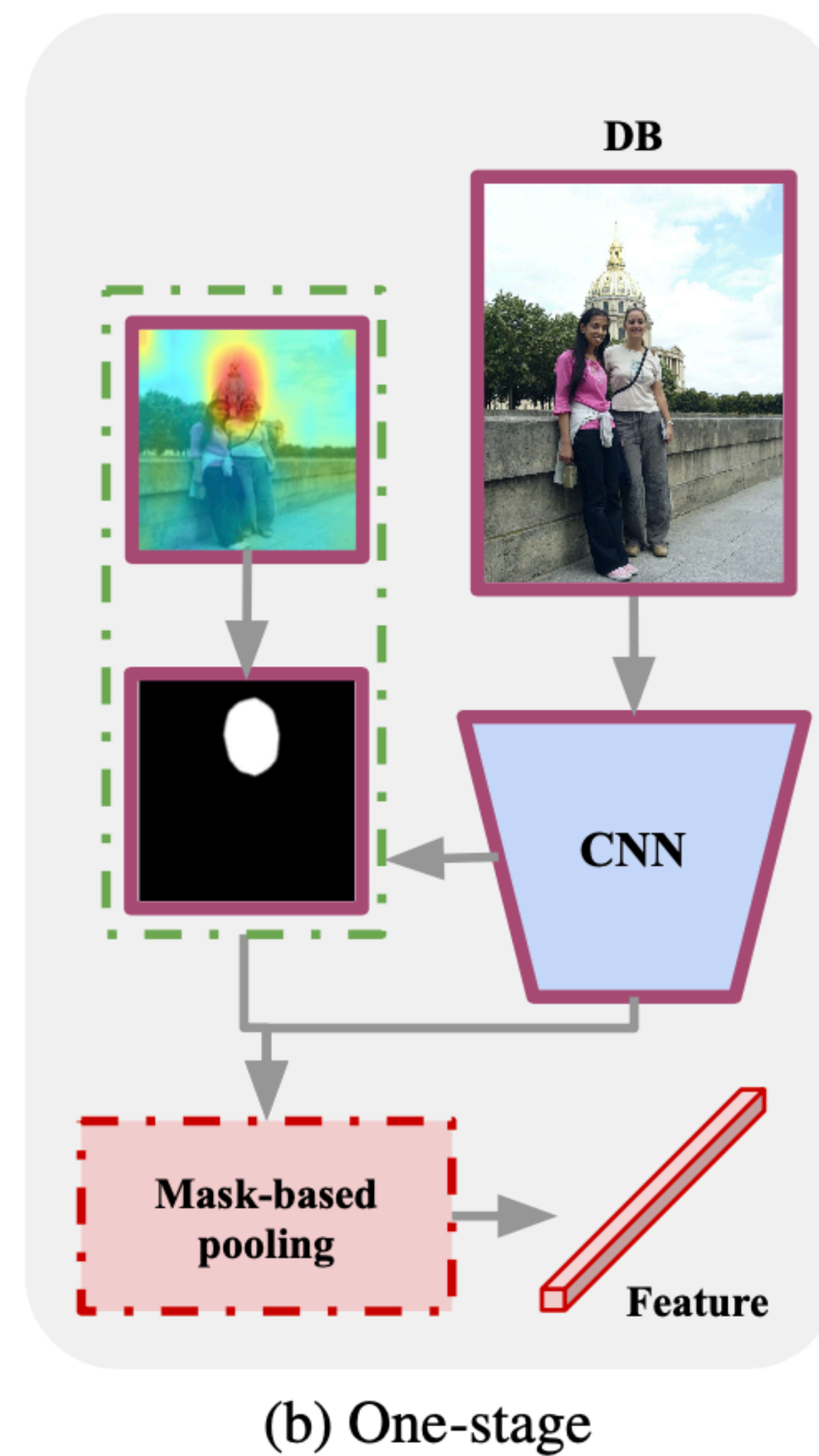
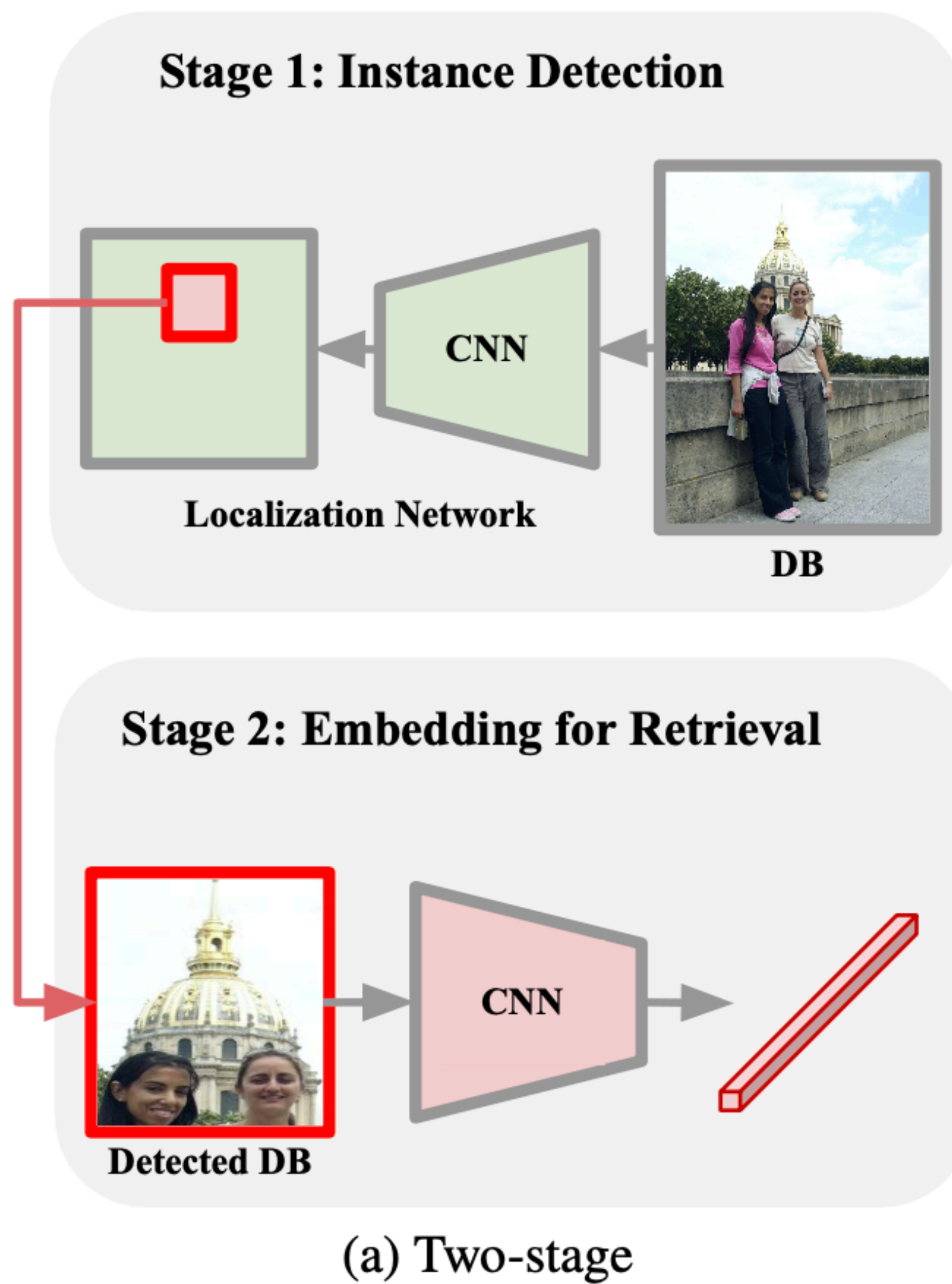
METHOD	TRAIN SET	BASE		MEDIUM				HARD				MEAN	DIFF
		Ox5k	Par6k	$\mathcal{R}Oxf$		$\mathcal{R}Par$		$\mathcal{R}Oxf$		$\mathcal{R}Par$			
		mAP	mAP	mAP	mP@10	mAP	mP@10	mAP	mP@10	mAP	mP@10		
Yokoo <i>et al.</i> [46]	GLDv2-clean	91.9	94.5	72.8	86.7	84.2	95.9	49.9	62.1	69.7	88.4	79.5	-5.4
Yokoo <i>et al.</i> [60] <sup>†</sup>	$\mathcal{R}GLDv2$ -clean	86.1	93.9	64.5	81.0	84.1	95.4	35.6	51.5	68.7	86.4	74.1	
SOLAR [58]	GLDv2-clean	–	–	79.7	–	88.6	–	60.0	–	75.3	–	75.9	-8
SOLAR [27] <sup>†</sup>	$\mathcal{R}GLDv2$ -clean	90.6	94.4	70.8	84.6	84.1	95.4	48.0	62.3	68.7	86.4	67.9	
GLAM [46]	GLDv2-clean	94.2	95.6	78.6	88.2	88.5	97.0	60.2	72.9	76.8	93.4	83.4	-4.1
GLAM [46] <sup>‡</sup>	$\mathcal{R}GLDv2$ -clean	90.9	94.1	72.2	84.7	83.0	95.0	49.6	61.6	65.6	87.6	79.3	
DOLG [47]	GLDv2-clean	–	–	78.8	–	87.8	–	58.0	–	74.1	–	74.7	-7.4
DOLG [59] <sup>†</sup>	$\mathcal{R}GLDv2$ -clean	88.3	93.9	70.8	85.3	83.2	95.4	47.4	60.0	67.9	87.4	67.3	
Token [58]	GLDv2-clean	–	–	82.3	–	75.6	–	66.6	–	78.6	–	75.8	-18.2
Token [58] <sup>†</sup>	$\mathcal{R}GLDv2$ -clean	84.3	90.0	61.4	76.4	75.8	94.0	36.9	55.2	54.4	81.0	57.6	

# Main Contributions 2



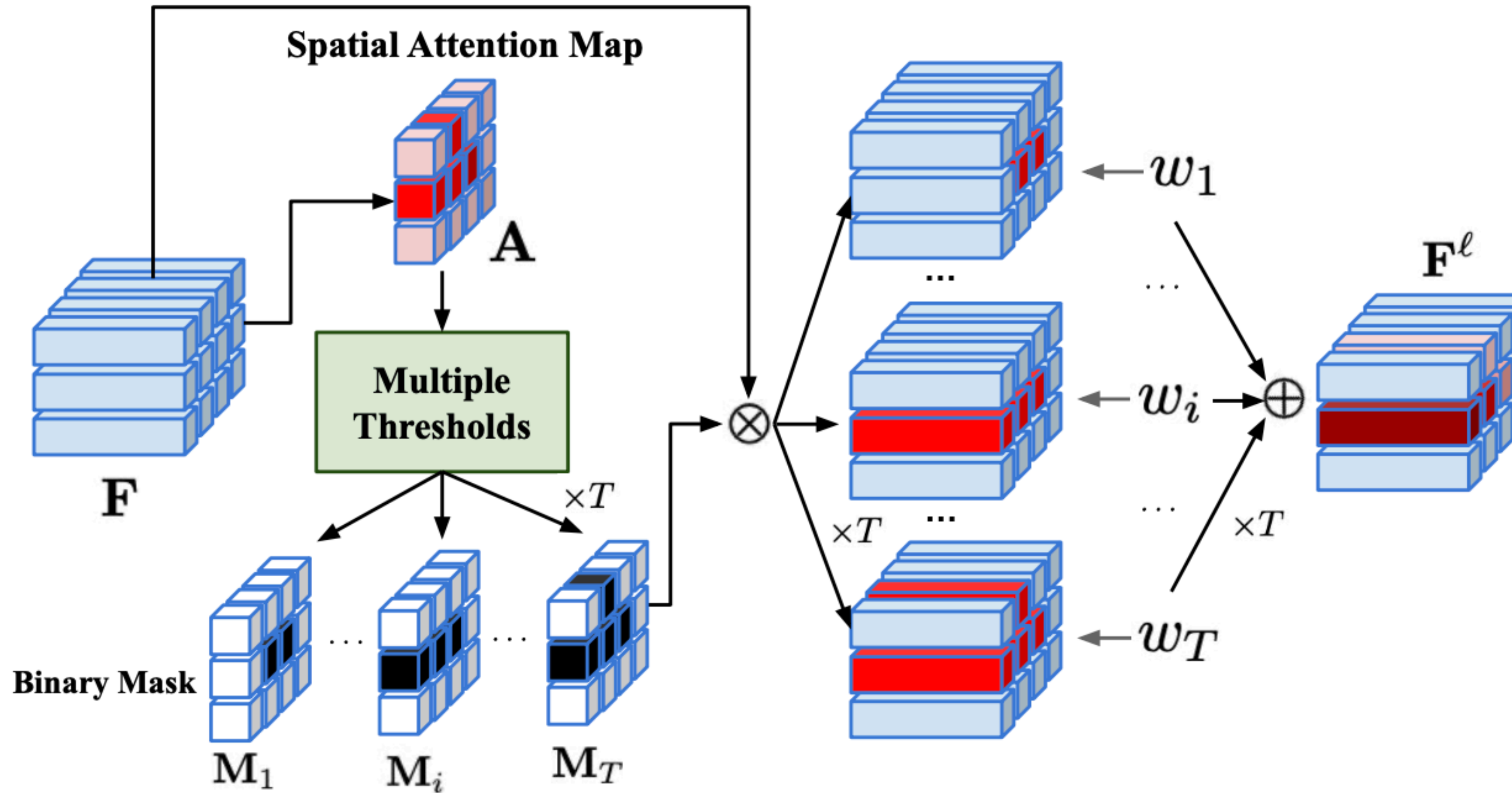


# Main Contributions 2



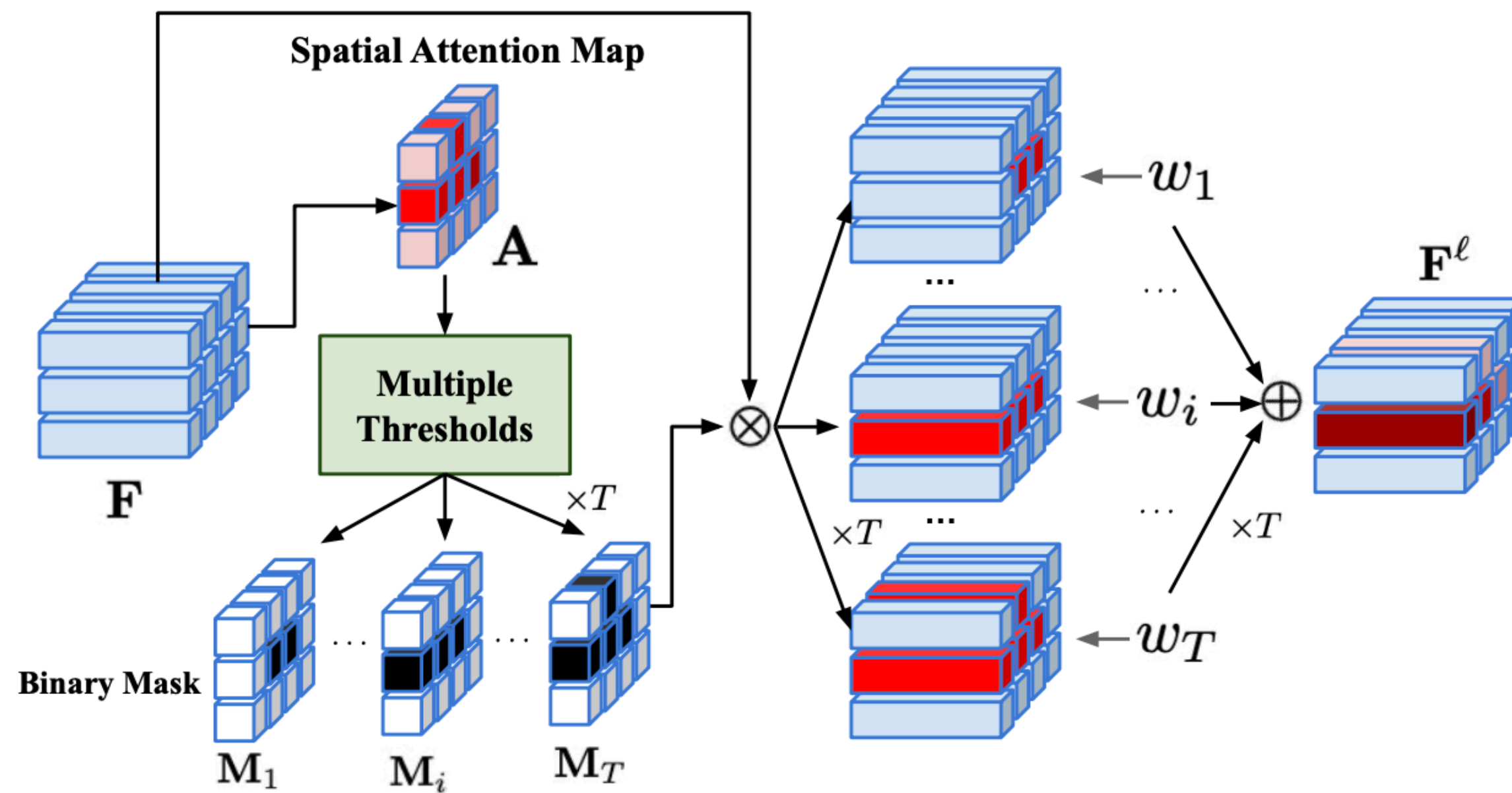
# Attentional Localization

- How spatial attention maps and masks are used



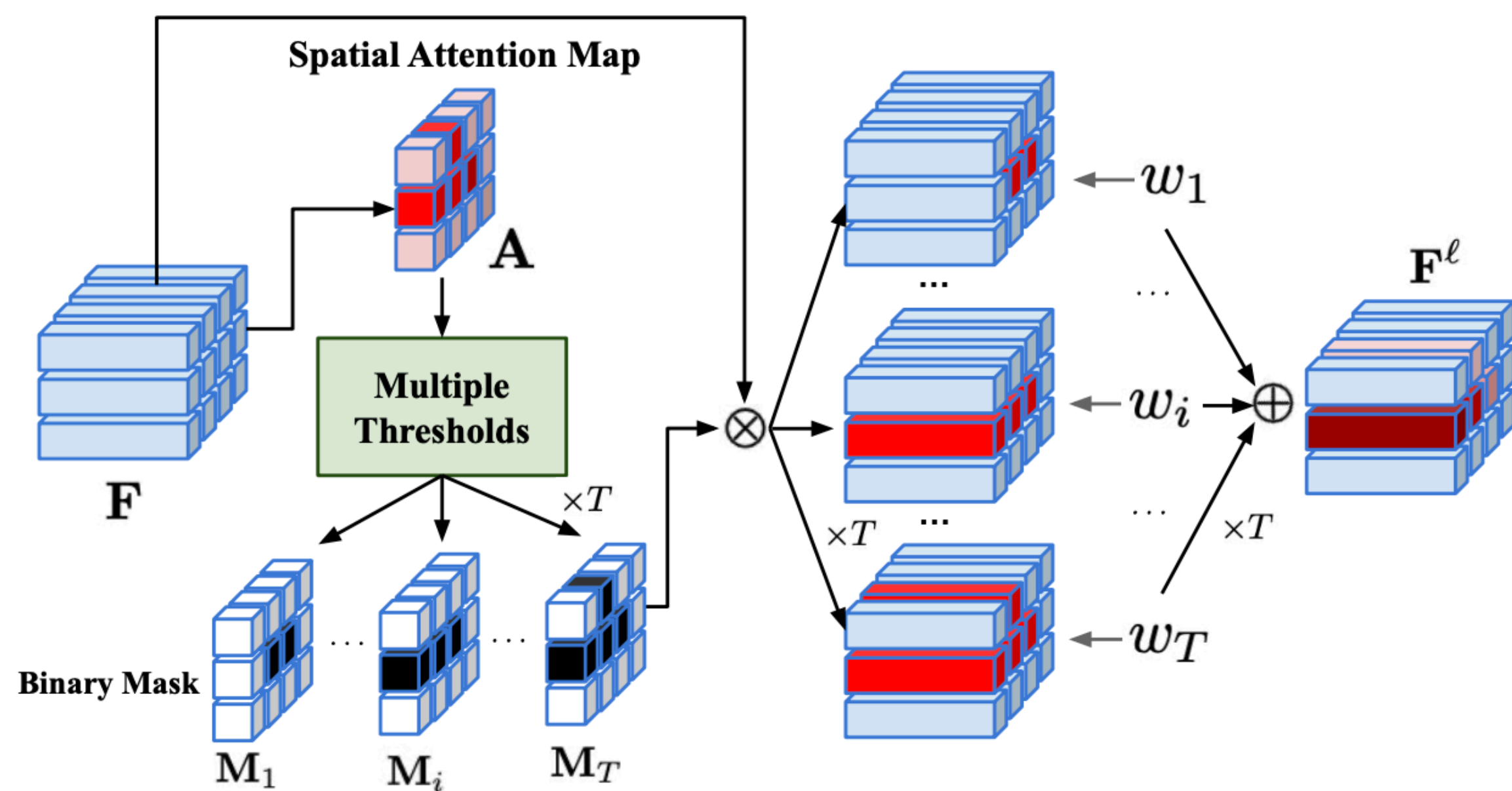


# Attentional Localization



$$A = \eta(\zeta(f^l(\mathbf{F}))) \in \mathbb{R}^{w \times h}$$

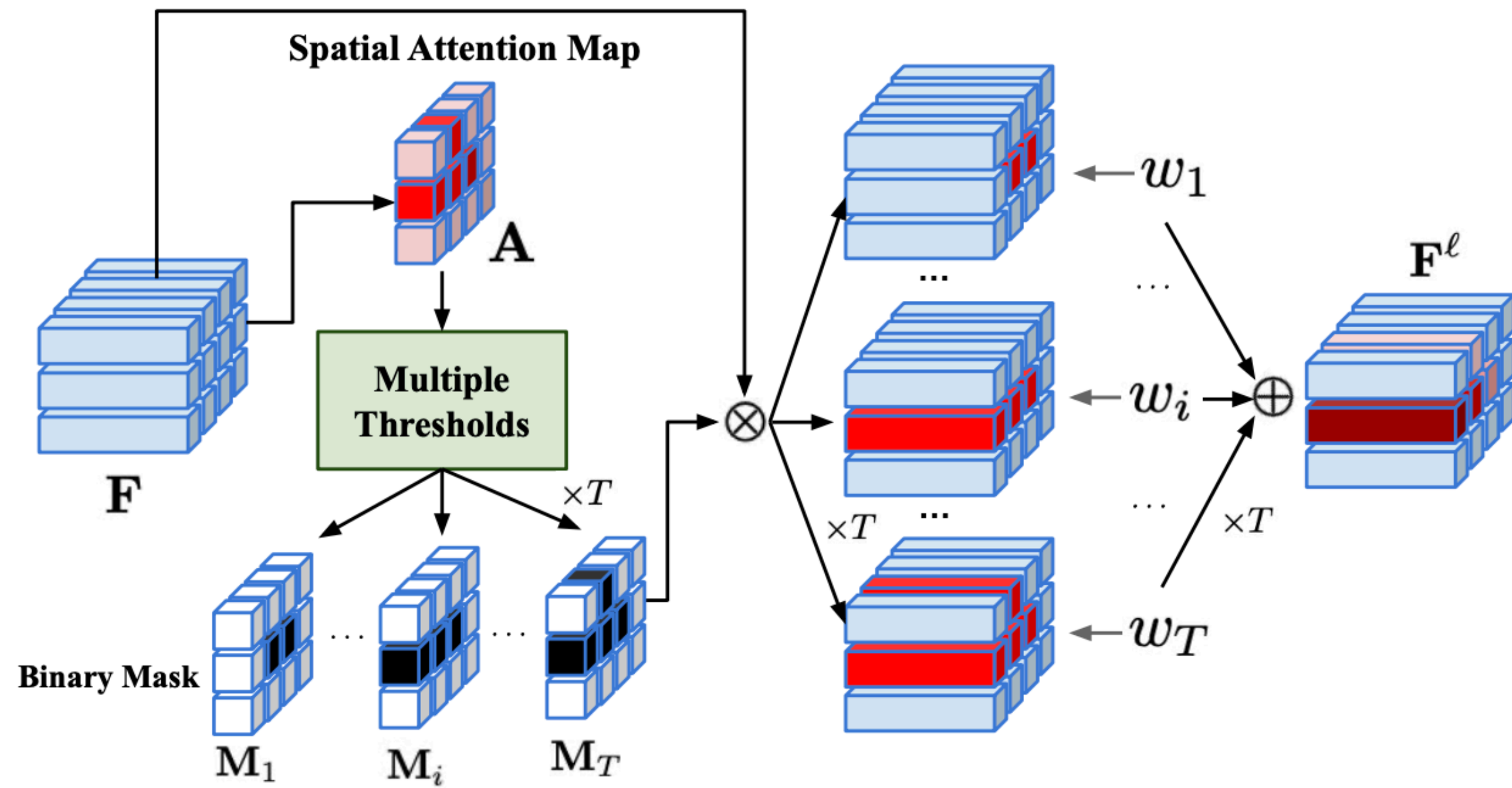
# Attentional Localization



$$\eta(X) := \frac{X - \min X}{\max X - \min X} \in \mathbb{R}^{w \times h}$$

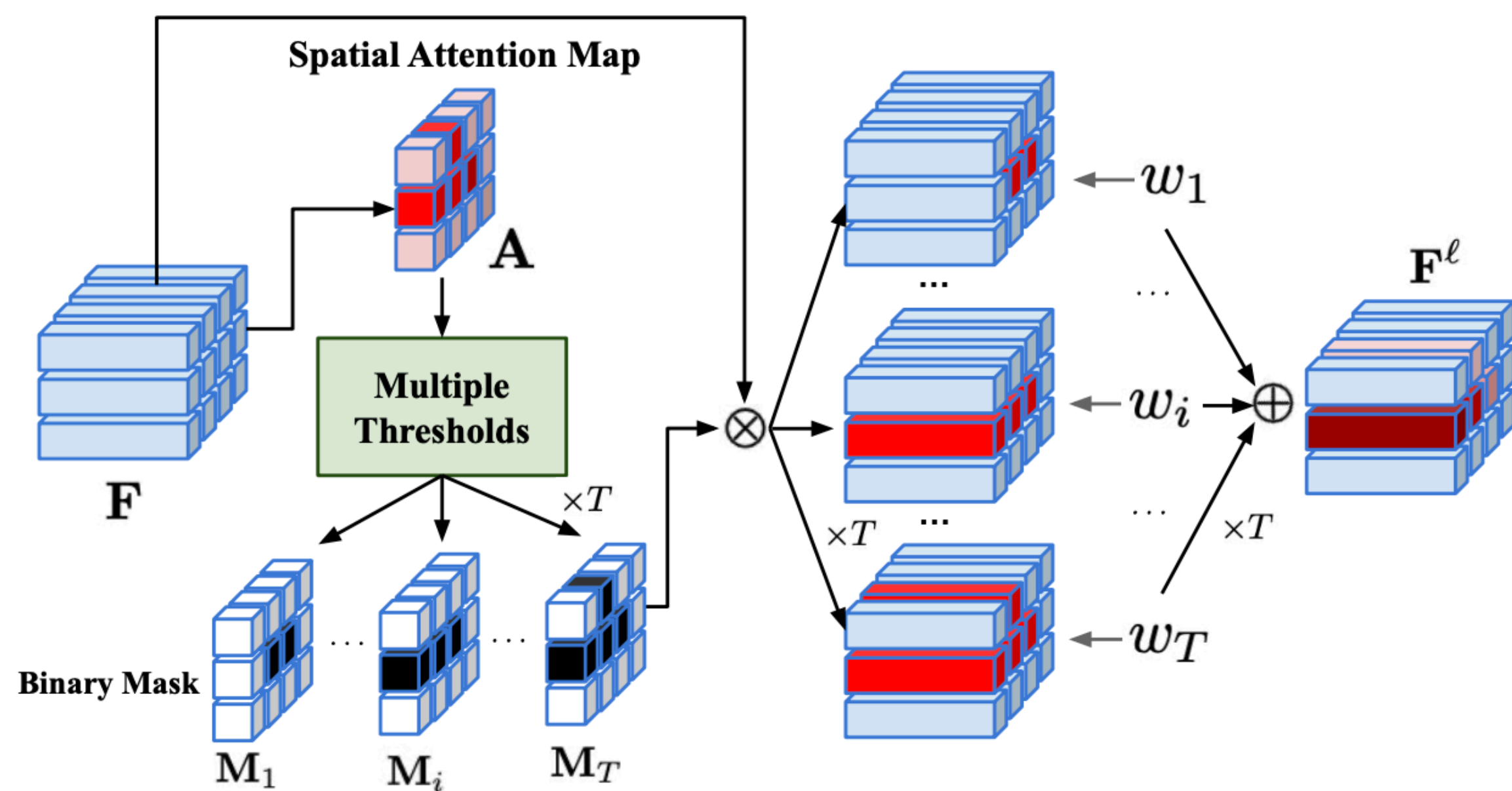


# Attentional Localization



$$M_i(\mathbf{p}) = \begin{cases} \beta, & \text{if } A(\mathbf{p}) < \tau_i \\ 1, & \text{otherwise} \end{cases}$$

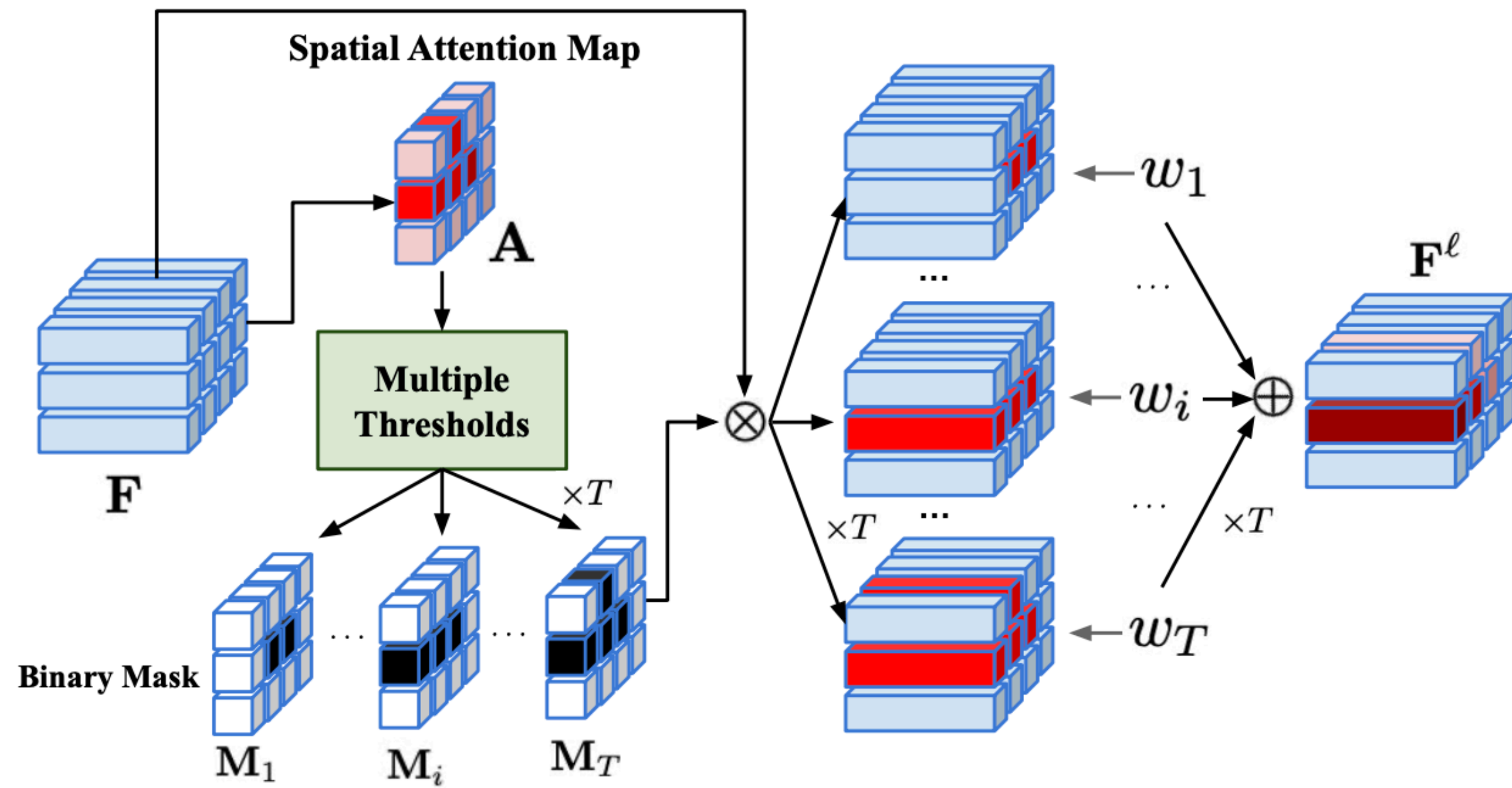
# Attentional Localization



$$\mathbf{F}^\ell = \mathcal{H}(M_1 \odot \mathbf{F}, \dots, M_T \odot \mathbf{F}) \in \mathbb{R}^{w \times h \times d}$$

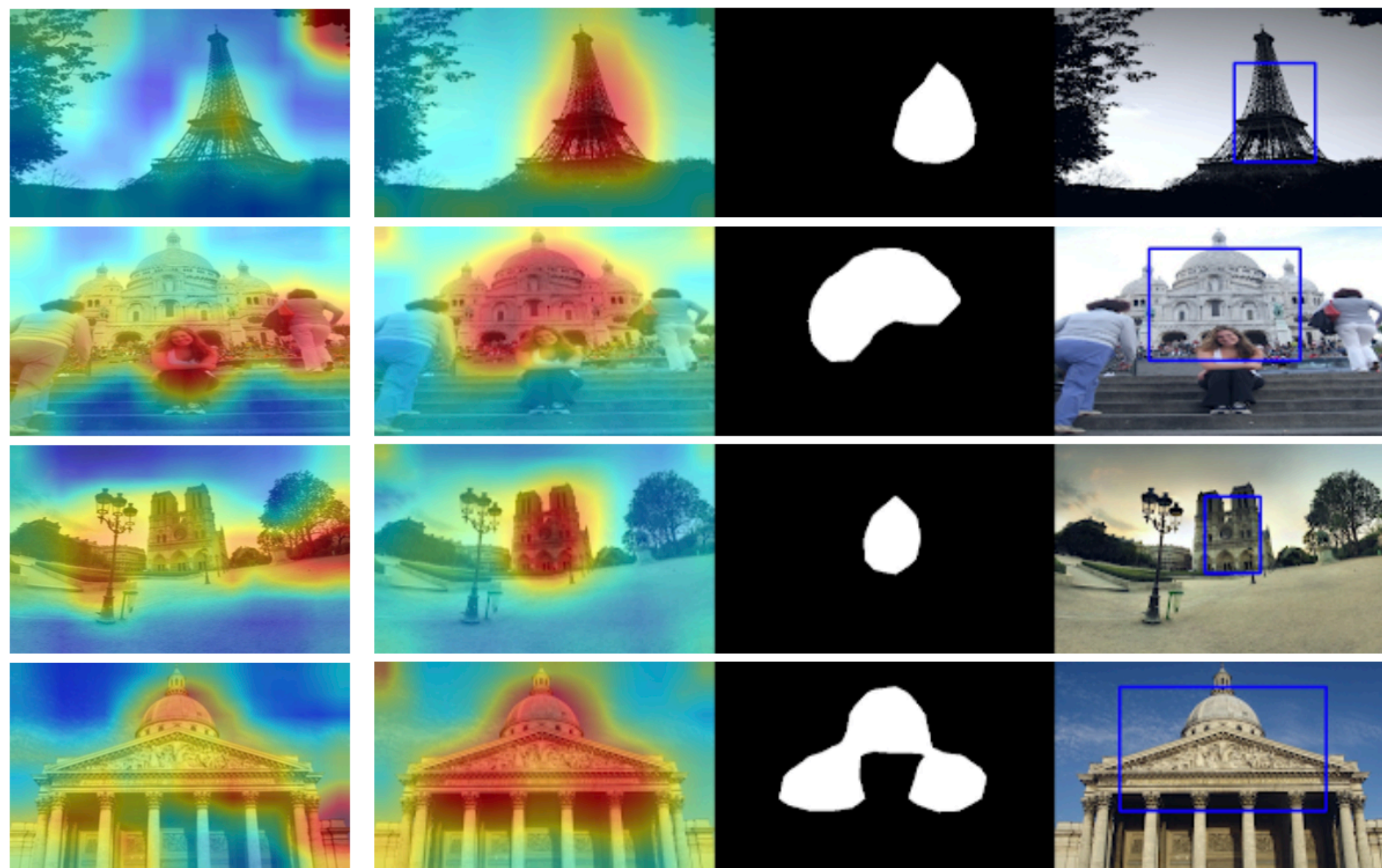


# Attentional Localization



$$\mathcal{H}(\mathbf{F}_1, \dots, \mathbf{F}_T) = \frac{w_1 \mathbf{F}_1 + \dots + w_T \mathbf{F}_T}{w_1 + \dots + w_T}$$

# Visual Comparison



(a)  $A$ , pre-trained

(b)  $A$ , ours

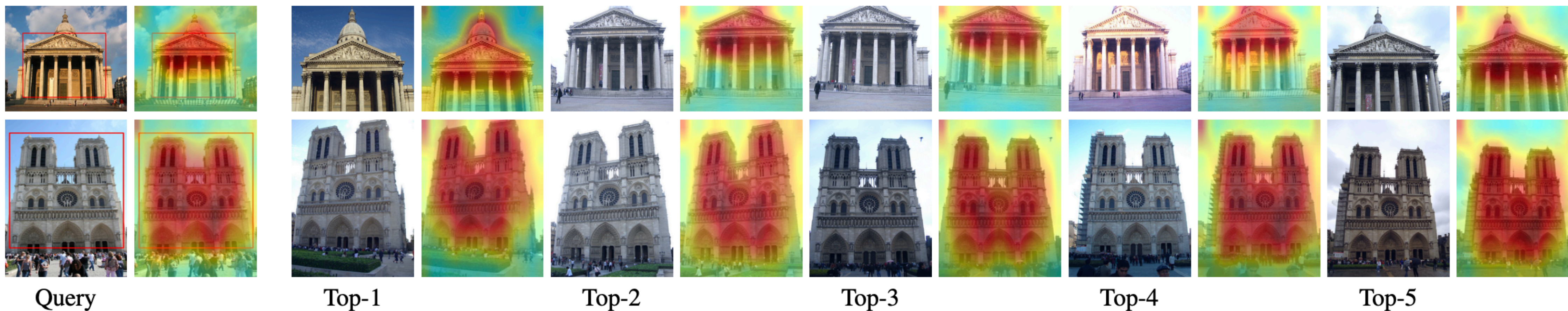
(c) Mask  $M_i$

(d) Bounding box



# Visualization

- Focus on relevant objects, ignoring background.





# Performance Comparison

- Superior performance on various datasets

METHOD	TRAIN SET	NET	POOLING	LOSS	FT	E2E	SELF	DIM	BASE		$\mathcal{R}_{\text{MEDIUM}}$		$\mathcal{R}_{\text{HARD}}$		MEAN	
									OXF5K	PAR6K	$\mathcal{R}_{\text{Oxf}}$	$\mathcal{R}_{\text{Par}}$	$\mathcal{R}_{\text{Oxf}}$	$\mathcal{R}_{\text{Par}}$		
LOCAL DESCRIPTORS																
HesAff-rSIFT-ASMK <sup>*</sup> +SP [34]	SfM-120k	R50	–	–	✓	–	–	–	–	–	60.6	61.4	36.7	35.0	–	
DELf-ASMK <sup>*</sup> +SP [34]	SfM-120k	R50	–	CLS	✓	–	–	–	–	–	<b>67.8</b>	<b>76.9</b>	<b>43.1</b>	<b>55.4</b>	–	
LOCAL DESCRIPTORS+D2R																
R-ASMK <sup>*</sup> [48]	NC-clean	R50	–	CLS,LOCAL	✓	–	–	–	–	–	69.9	<b>78.7</b>	45.6	<b>57.7</b>	–	
R-ASMK <sup>*</sup> +SP [48]	NC-clean	R50	–	CLS,LOCAL	✓	–	–	–	–	–	<b>71.9</b>	78.0	<b>48.5</b>	54.0	–	
GLOBAL DESCRIPTORS																
DIR [47]	SfM-120k	R101	RMAC	TP	✓	–	–	2048	79.0	86.3	53.5	68.3	25.5	42.4	59.2	
Radenovic <i>et al.</i> [36, 34]	SfM-120k	R101	GeM	SIA	–	–	–	2048	87.8	<b>92.7</b>	64.7	77.2	38.5	56.3	69.5	
AGeM [9]	SfM-120k	R101	GeM	SIA	–	–	–	2048	–	–	<b>67.0</b>	<b>78.1</b>	<b>40.7</b>	<b>57.3</b>	–	
SOLAR [47]	SfM-120k	R101	GeM	TP,SOS	✓	–	–	2048	78.5	86.3	52.5	70.9	27.1	46.7	60.3	
GLAM [46]	SfM-120k	R101	GeM	AF	–	–	–	512	<b>89.7</b>	91.1	66.2	77.5	39.5	54.3	<b>69.7</b>	
DOLG [47]	SfM-120k	R101	GeM,GAP	AF	–	–	–	512	72.8	74.5	46.4	56.6	18.1	26.6	49.2	
GLOBAL DESCRIPTORS+D2R																
Mei <i>et al.</i> [26]	[O]	R101	FC	CLS	–	–	–	4096	38.4	–	–	–	–	–	–	
Salvador <i>et al.</i> [43]	Pascal VOC	V16	GSP	CLS,LOCAL	–	✓	–	512	67.9	72.9	–	–	–	–	–	
Chen <i>et al.</i> [4]	OpenImageV4 [17]	R50	MAC	MSE	–	✓	–	2048	50.2	65.2	–	–	–	–	–	
Liao <i>et al.</i> [22]	Oxford,Paris	A,V16	CroW	CLS,LOCAL	–	–	–	768	80.1	90.3	–	–	–	–	–	
DIR+RPN [8]	NC-clean	R101	RMAC	TP	✓	–	–	2048	<b>85.2</b>	<b>94.0</b>	–	–	–	–	–	
<b>CiDeR (Ours)</b>	SfM-120k	R101	GeM	AF	–	✓	✓	2048	<b>89.9</b>	92.0	<b>67.3</b>	<b>79.4</b>	<b>42.4</b>	57.5	<b>71.4</b>	
<b>CiDeR-FT (Ours)</b>	SfM-120k	R101	GeM	AF	✓	✓	✓	2048	<b>92.6</b>	<b>95.1</b>	<b>76.2</b>	<b>84.5</b>	<b>58.9</b>	<b>68.9</b>	<b>79.4</b>	

# Performance Comparison

- Superior performance on various datasets

METHOD	BASE		MEDIUM								HARD							
	Ox5k	Par6k	$\mathcal{R}Oxf$		$\mathcal{R}Oxf + \mathcal{R}1M$		$\mathcal{R}Par$		$\mathcal{R}Par + \mathcal{R}1M$		$\mathcal{R}Oxf$		$\mathcal{R}Oxf + \mathcal{R}1M$		$\mathcal{R}Par$		$\mathcal{R}Par + \mathcal{R}1M$	
	mAP	mAP	mAP	mP@10	mAP	mP@10	mAP	mP@10	mAP	mP@10	mAP	mP@10	mAP	mP@10	mAP	mP@10	mAP	mP@10
GLOBAL DESCRIPTORS (SFM-120K)																		
DIR [47]	79.0	86.3	53.5	76.9	–	–	68.3	97.7	–	–	25.5	42.0	–	–	42.4	83.6	–	–
Filip <i>et al.</i> [36, 34]	87.8	92.7	64.7	<b>84.7</b>	<b>45.2</b>	<b>71.7</b>	77.2	<b>98.1</b>	<b>52.3</b>	<b>95.3</b>	38.5	<b>53.0</b>	<b>19.9</b>	<b>34.9</b>	56.3	<b>89.1</b>	24.7	<b>73.3</b>
AGeM [9]	–	–	<b>67.0</b>	–	–	–	<b>78.1</b>	–	–	–	<b>40.7</b>	–	–	–	57.3	–	–	–
SOLAR [47]	78.5	86.3	52.5	73.6	–	–	70.9	<b>98.1</b>	–	–	27.1	41.4	–	–	46.7	83.6	–	–
GeM [47]	79.0	82.6	54.0	72.5	–	–	64.3	92.6	–	–	25.8	42.2	–	–	36.6	67.6	–	–
GLAM [47]	<b>89.7</b>	91.1	66.2	–	–	–	77.5	–	–	–	39.5	–	–	–	54.3	–	–	–
DOLG [47]	72.8	74.5	46.4	66.8	–	–	56.6	91.1	–	–	18.1	27.9	–	–	26.6	62.6	–	–
<b>CiDeR (Ours)</b>	<b>89.9</b>	92.0	<b>67.3</b>	<b>85.1</b>	<b>50.3</b>	<b>75.5</b>	<b>79.4</b>	97.9	51.4	<b>95.7</b>	<b>42.4</b>	<b>56.4</b>	<b>22.4</b>	<b>35.9</b>	57.5	87.1	22.4	69.4
<b>CiDeR-FT (Ours)</b>	<b>92.6</b>	<b>95.1</b>	<b>76.2</b>	<b>87.3</b>	<b>60.5</b>	<b>78.6</b>	<b>84.5</b>	98.0	<b>56.9</b>	<b>95.9</b>	<b>58.9</b>	<b>71.1</b>	<b>36.8</b>	<b>55.7</b>	<b>68.9</b>	<b>91.3</b>	<b>30.1</b>	<b>73.9</b>
GLOBAL DESCRIPTORS ( $\mathcal{R}GLDV2$ -CLEAN)																		
Yokoo <i>et al.</i> [60] <sup>†</sup> (Base)	86.1	93.9	64.5	81.0	51.3	72.1	84.1	<b>95.4</b>	54.2	90.3	35.6	51.5	22.2	42.9	<b>68.7</b>	86.4	27.4	66.9
SOLAR [27] <sup>†</sup>	90.6	<b>94.4</b>	70.8	84.6	55.8	76.1	80.3	94.6	57.6	<b>92.0</b>	48.0	<b>62.3</b>	30.3	45.3	61.8	83.9	30.7	71.6
GLAM [46] <sup>‡</sup>	<b>90.9</b>	94.1	<b>72.2</b>	84.7	<b>58.6</b>	76.1	83.0	95.0	<b>58.6</b>	91.7	<b>49.6</b>	61.6	<b>34.1</b>	<b>50.9</b>	65.6	<b>87.6</b>	<b>33.3</b>	72.1
DOLG [59] <sup>†</sup>	88.3	93.9	70.8	<b>85.3</b>	57.3	<b>76.8</b>	<b>83.2</b>	<b>95.4</b>	57.3	<b>92.0</b>	47.4	60.0	29.5	46.2	67.9	87.4	32.7	<b>72.4</b>
Token [58] <sup>†</sup>	81.2	89.6	60.8	77.7	44.0	60.9	75.8	94.3	44.1	86.9	37.3	54.1	23.2	37.7	54.8	81.3	19.7	54.4
<b>CiDeR (Ours)</b>	89.8	<b>94.6</b>	<b>73.7</b>	<b>85.5</b>	<b>58.6</b>	76.3	<b>84.6</b>	<b>96.7</b>	<b>59.0</b>	<b>95.1</b>	<b>54.9</b>	<b>66.6</b>	<b>34.6</b>	<b>54.7</b>	<b>68.5</b>	<b>89.1</b>	<b>33.5</b>	<b>76.9</b>
<b>CiDeR-FT (Ours)</b>	<b>90.9</b>	<b>96.1</b>	<b>77.8</b>	<b>88.0</b>	<b>61.8</b>	<b>78.0</b>	<b>87.4</b>	<b>97.0</b>	<b>61.6</b>	<b>94.3</b>	<b>61.9</b>	<b>70.4</b>	<b>39.4</b>	<b>56.8</b>	<b>75.3</b>	<b>90.0</b>	<b>35.8</b>	<b>72.7</b>



# Conclusion