

无人机多目标追踪、多智能体 协同感知相关工作介绍

王永才

中国人民大学 信息学院 计算机系

ycw@ruc.edu.cn

中国人民大学计算机系介绍

信息学院

计算机科学与技术系

经济信息管理系

数据工程与知识工程教育部重点实验室

高瓴人工智能学院

大数据管理与分析方法研究北京市重点实验室

新一代智能搜索与推荐教育部工程研究中心

北京市科技进步一等奖、国家科技进步二等奖、电子部科技进步特等奖、北京市科技进步二等奖（2次）、教育部科技进步二等奖（2次），建国70周年纪念奖。



INLAB: 智能网络与优化实验室介绍



李德英教授



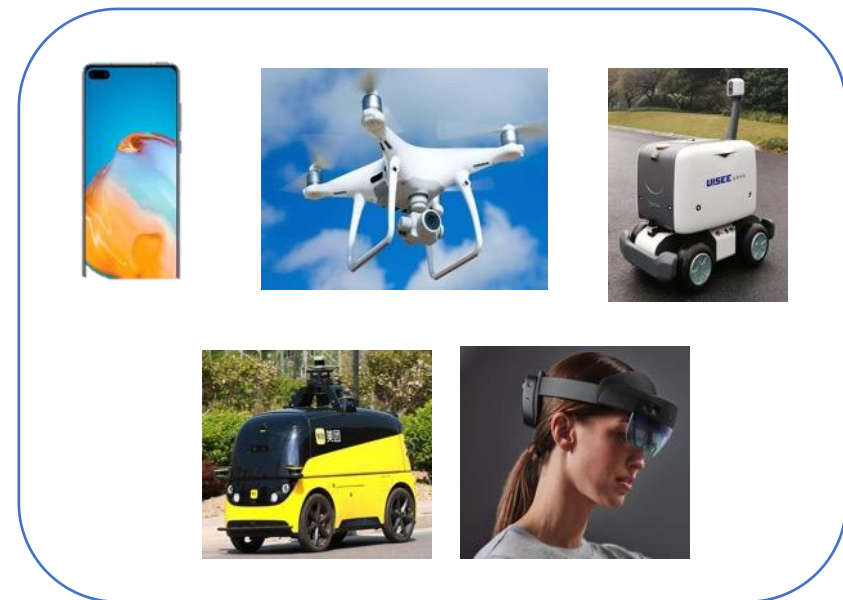
王永才副教授



周春来副教授



陈文萍讲师



主要研究方向：新型AIOT智能感知网络
智能感知、多机协同感知等

• 近期承担的主要项目

1. 国家自然科学基金重点基金，面向大数据机器学习的不确定性建模理论与方法，子课题，2018-2022
2. 国家自然科学基金面上项目，社会网络影响力传播机制及最优化方法研究，面上项目，2021-2024
3. 国家自然科学基金面上项目，基于模块拼接的群智SLAM关键问题研究，2020-2023
4. 国家科技支撑计划子课题，智能船自动驾驶辅助系统，2020-2022
5. 国家自然科学基金面上项目，大规模社会网络信息处理优化技术研究，2017-2020
6. 国家自然科学基金面上项目，紧耦合多源同步定位与地图构建的融合优化关键问题研究，2017-2020

个人简介



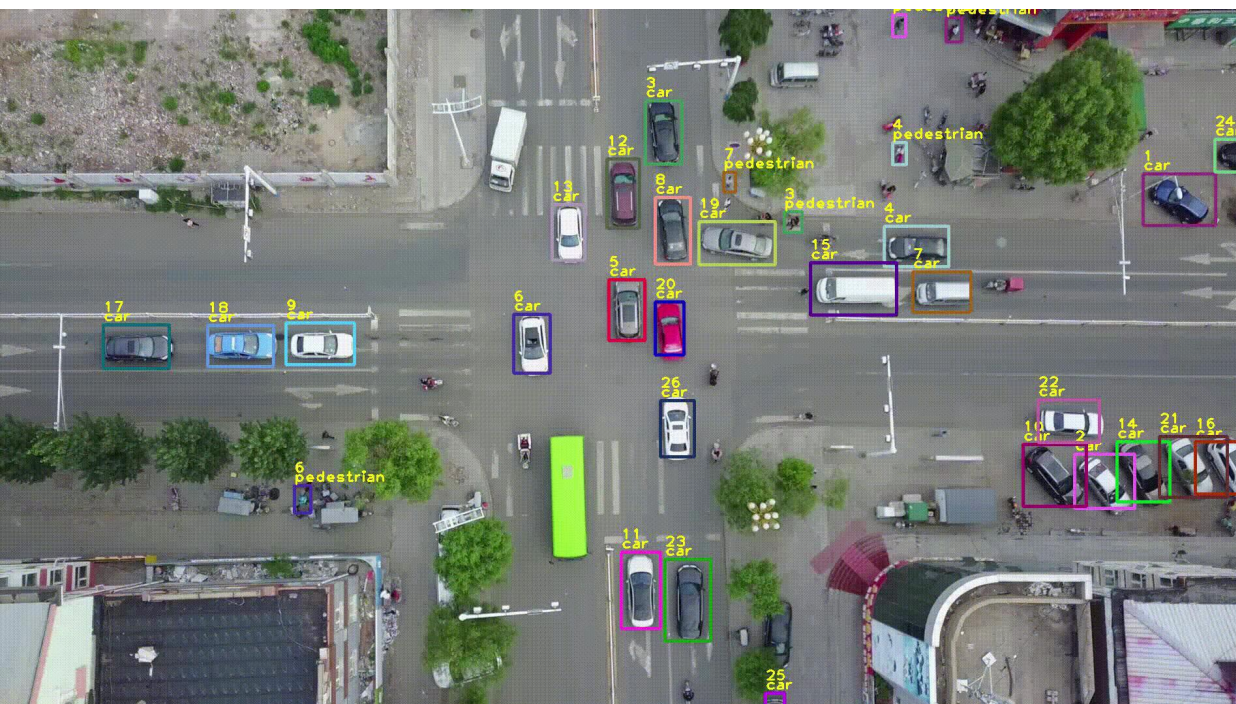
1997.08–2001.07	清华大学自动化系，本科
2001.09–2006.12	清华大学自动化系，博士生
2007.01–2009.08	NEC中国研究院，副研究员
2009.08–2015.08	清华大学交叉信息研究院，助理研究员
2015.08–2022.06	中国人民大学信息学院，副教授
2014.01–2014.08	美国康奈尔大学，访问学者

主要研究领域为**多智能体协同感知、图优化、SLAM系统**等，在国内外知名期刊和会议**发表论文100余篇**，已授权专利10余项。研究成果被应用于**智能车、智能船领域**。主持多项国家自然科学基金面上项目，国家科技支撑计划子课题，2021年获得交通运输部航海学会**技术发明奖一等奖**，2022年获得交通运输部航海学会**科技进步二等奖**。Email: ycw@ruc.edu.cn，**电话微信：18910215881**

目录

- DroneMOT: 无人机多目标追踪系统
- GSLAMOT: 同步定位建图与多目标追踪系统
- SAHNet/RoCo: 多车协同感知系统
- MobiSketch: 基于点线面融合的语义3D建图
- CoISLAM: 多手机协同SLAM系统

可视化对比: DroneMOT and UAVMOT



DroneMOT 我们的方法

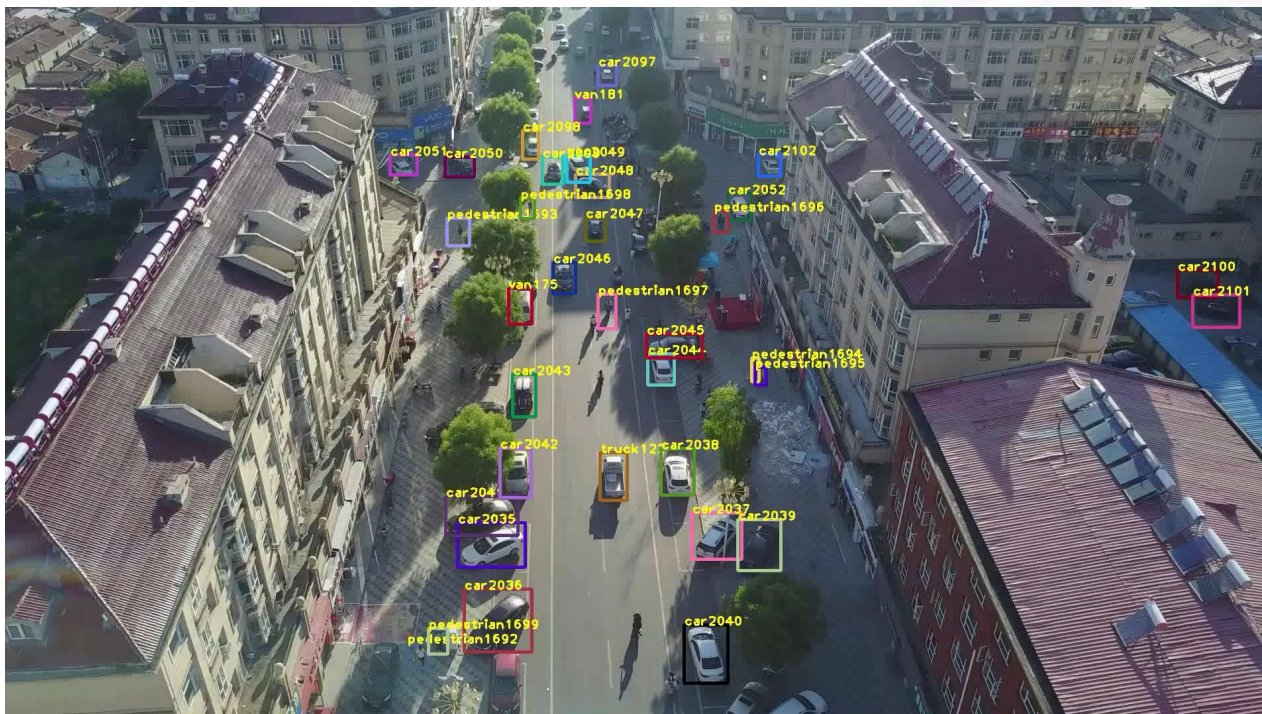
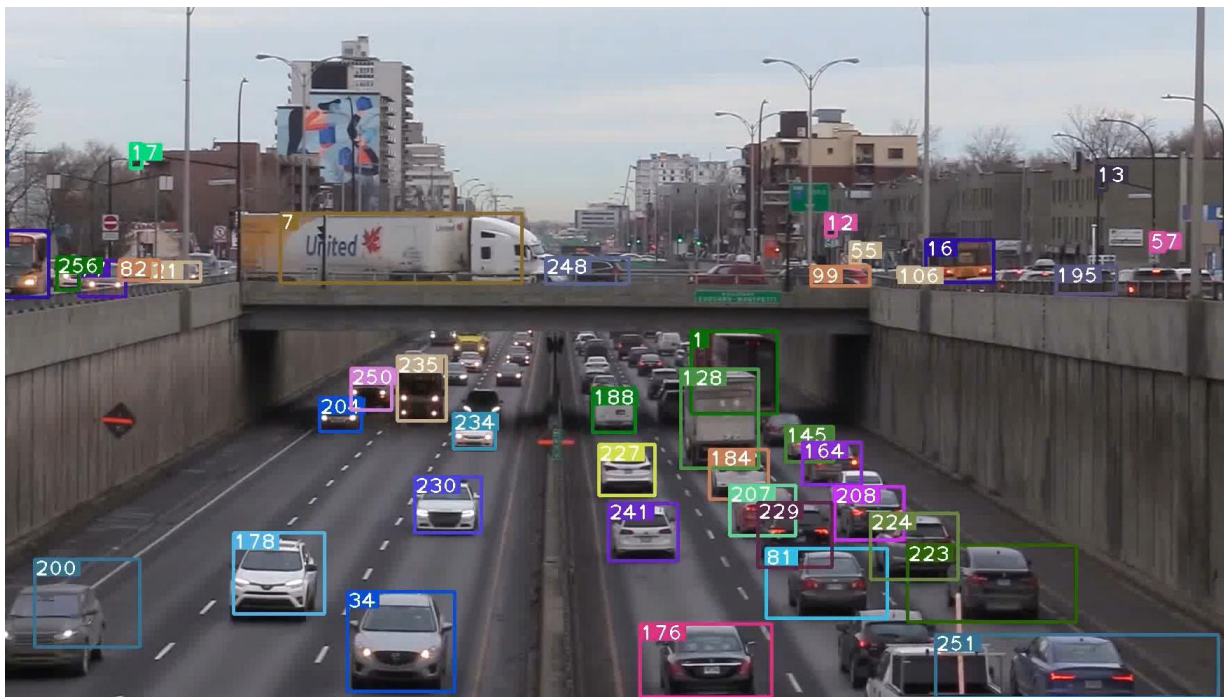
无人机旋转 情况下跟踪的更为准确



UAVMOT

无人机旋转

Multi-Object Tracking (MOT) : 多目标追踪的应用

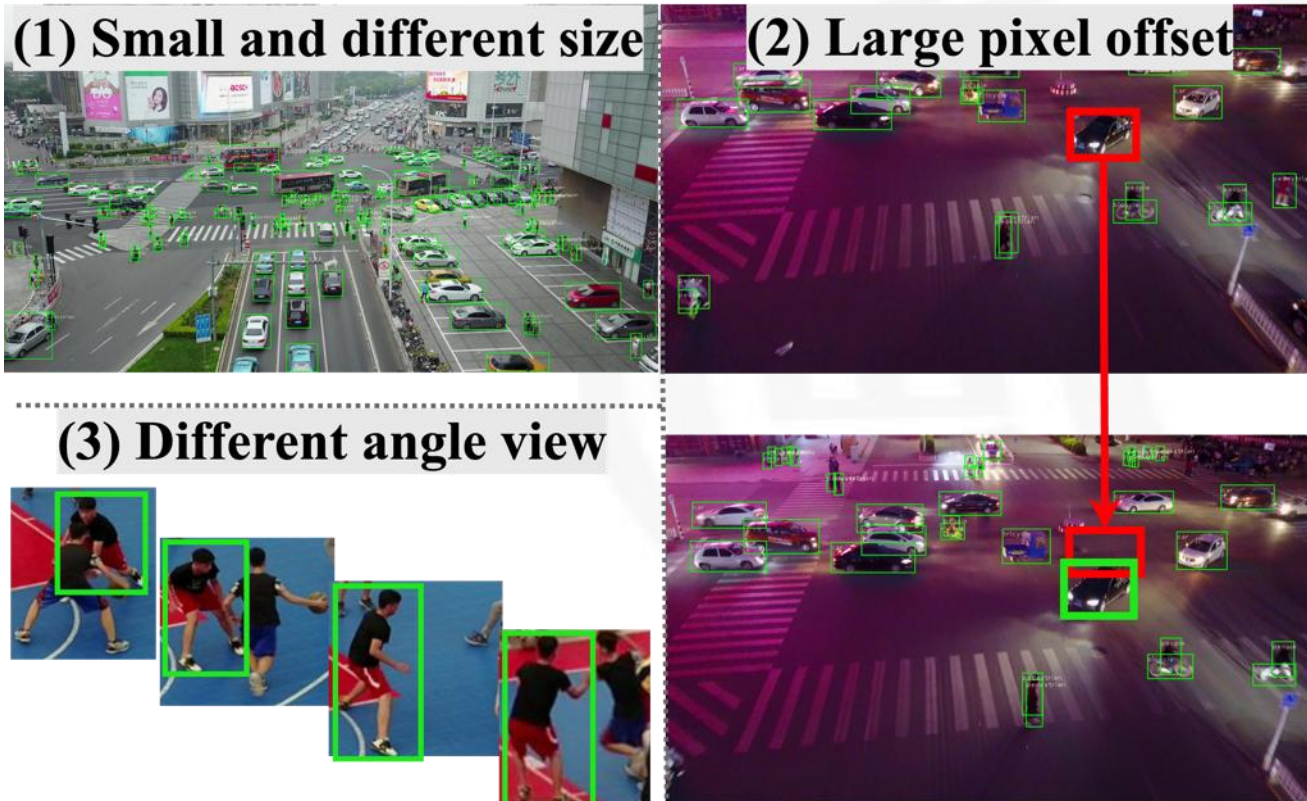


既要检测移动目标，又要在连续画面中关联和追踪移动目标

主要用于车辆追踪，行人追踪等

在无人机(Drones)上进行多目标追踪的**主要挑战**

- I. 目标像素尺寸小，**检测困难**
- II. 无人机运动不规则，**目标轨迹预测基于Kalman滤波等预测不准确。**
- III. 由于无人机移动，**目标不间断的进出无人机视野。**



现有的无人机上的多目标跟踪方法

无人机多目标跟踪算法主要基于监控相机算法，遵循**基于检测的跟踪**范式对**目标检测**和**特征匹配**环节进行了特别的改进。

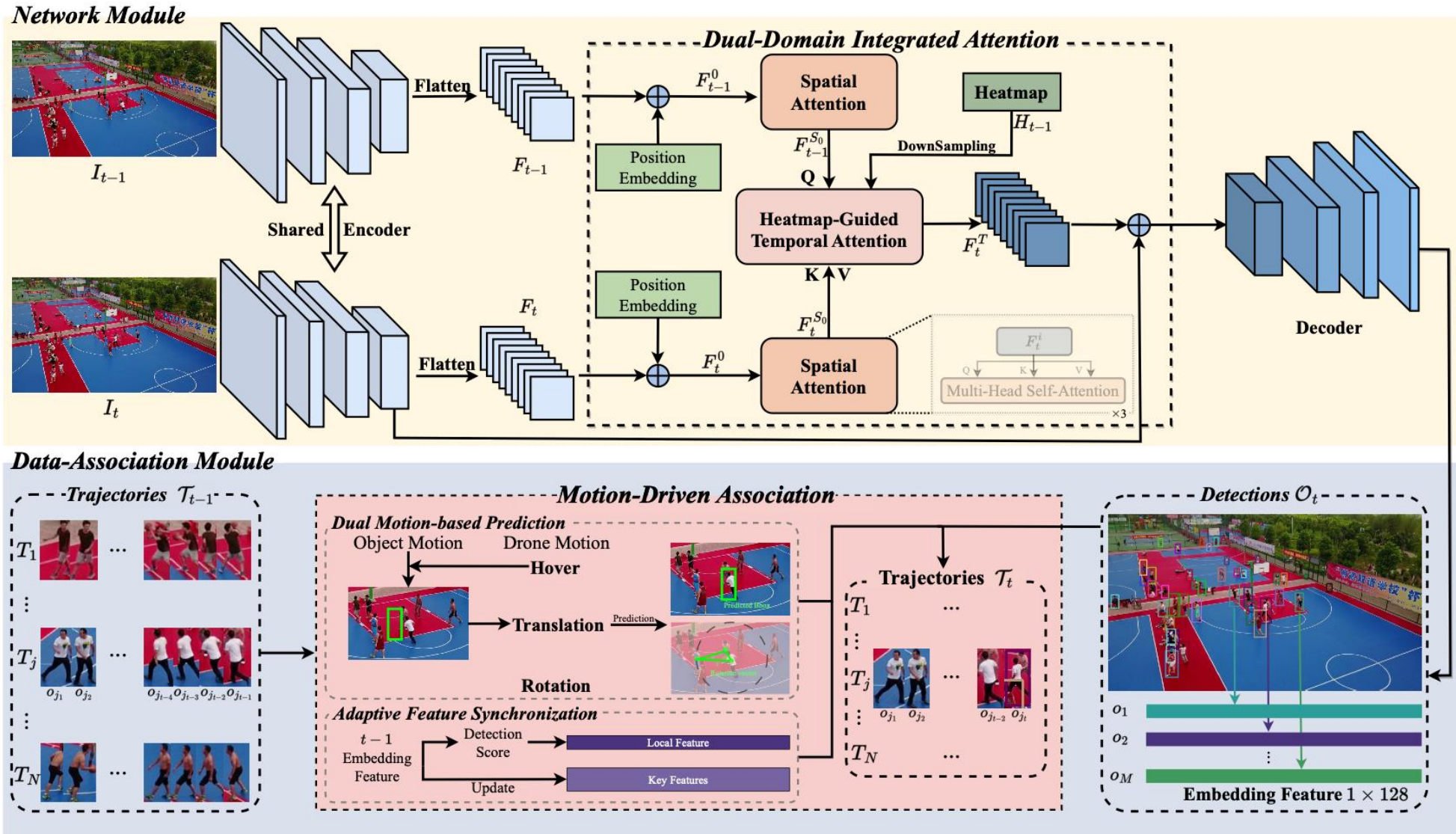
- Luo等人引入了**时空注意力模块的残差网络作为特征提取网络**，将YOLOv5与DeepSORT结合，加强了对感知小目标检测和外观特征提取的能力。
- FOLT采纳了一种轻量级**光流提取网络**，以最小的成本提取物体检测特征和运动特征。在神经网络中使用**光流引导特征增强**被设计为基于其光流来增强对象检测特征，从而提高了小对象的检测。然后还提出了**光流引导运动预测**来预测目标在下一帧中的位置，从而提高了相邻帧之间位移非常大的目标的跟踪性能。
- PAS Tracker通过融入一个**额外的ReID网络**来提取物体的特征，并综合考虑**位置、外观和尺寸**信息来表示物体，从而增强目标匹配的精度。
- GLOA则设计了一个全局-局部感知检测器，旨在捕获输入帧中的尺度变化特征信息，使用专门设计的全局局部感知块（CNN和Transformer结构）从输入帧中提取**不同尺度特征信息**。然后，它通过向预测头添加身份嵌入分支来输出更具辨别力的身份信息，从而更好地处理尺度方差问题并提高对象跟踪精度。

无人机上的多目标跟踪

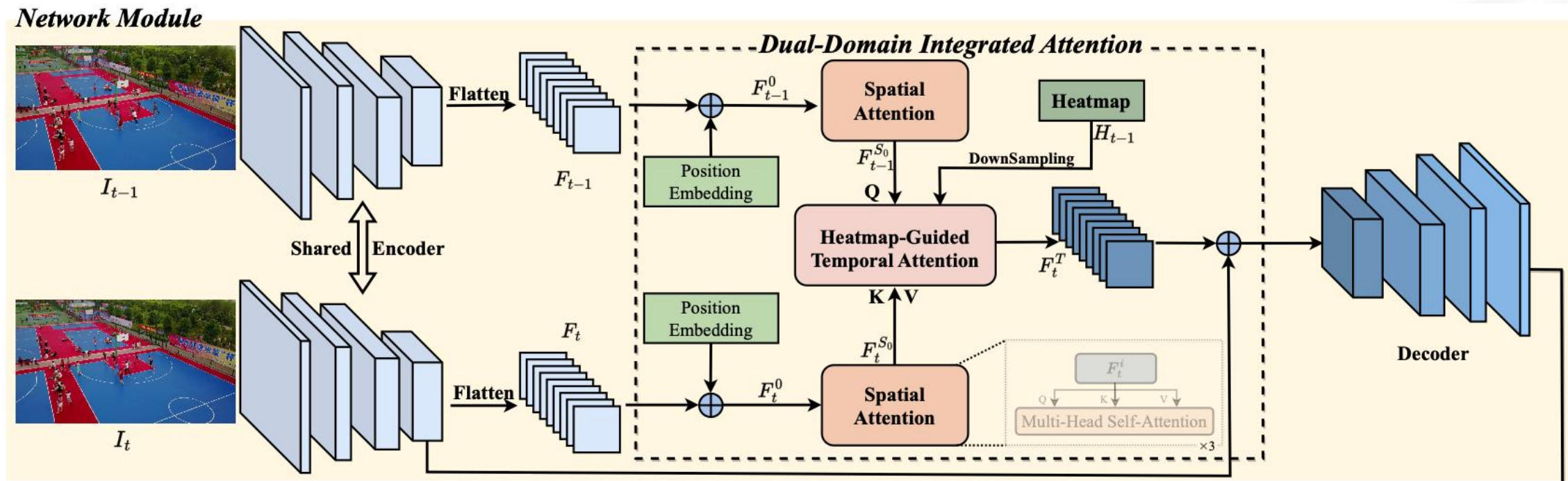
大多数算法关注于如何设计强大的目标检测网络，而对**特征匹配环节**研究不够充分。

- TrackletNet Tracker利用**时间和外观**信息，使用**图神经网络**来跟踪无人机检测到的目标，然后根据多视图立体技术估计的组平面来定位跟踪的地面物体，展现出在复杂环境下进行多目标跟踪的有效性。
- UAVMOT通过利用**两个相邻帧之间的关联层**来加强目标特征，引入了ID特征更新模块来增强对象的特征关联，同时使用梯度平衡焦点损失来解决**不平衡类别和小物体检测**问题。在匹配阶段，为了更好地处理无人机视图下的复杂运动，并开发了一个**自适应运动滤波器**，以实现精确的物体ID匹配。

我们提出的 DroneMOT 系统框架



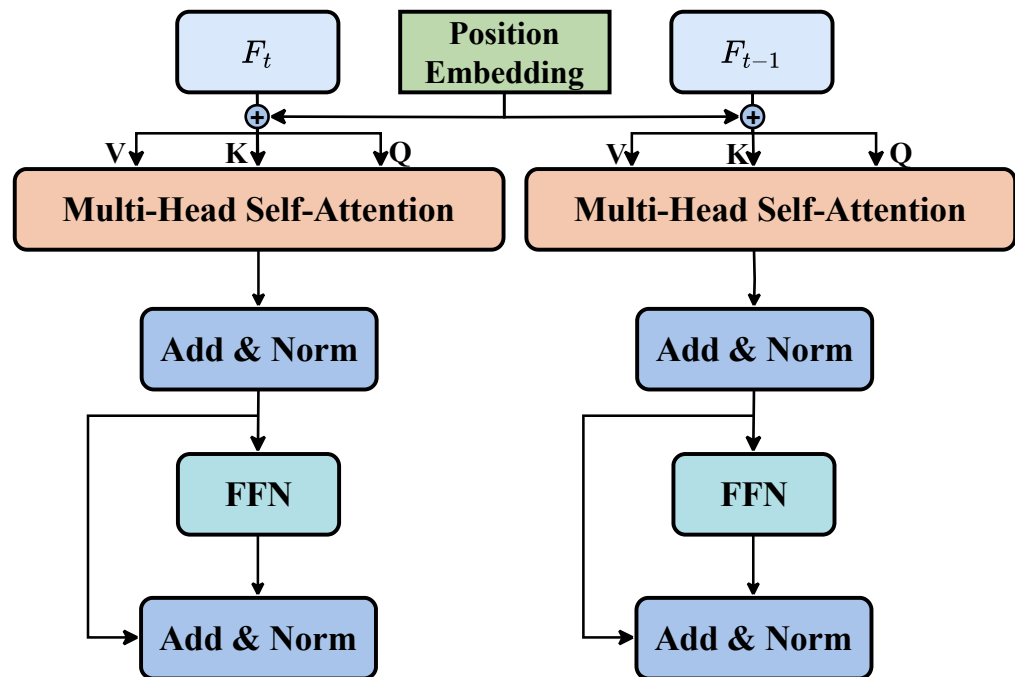
双语注意力模块



在网络模块中，我们使用DLA34网络作为骨干，它是一个**编码器-解码器架构**。

- 输入：来自 **t 帧和 t-1 帧的图像 I_t 和 $I_{(t-1)}$** 。
- 编码器使用共享的卷积层从图像 $\{I_{(t-1)}, I_t\}$ 中**提取局部特征**。

空间注意力模块



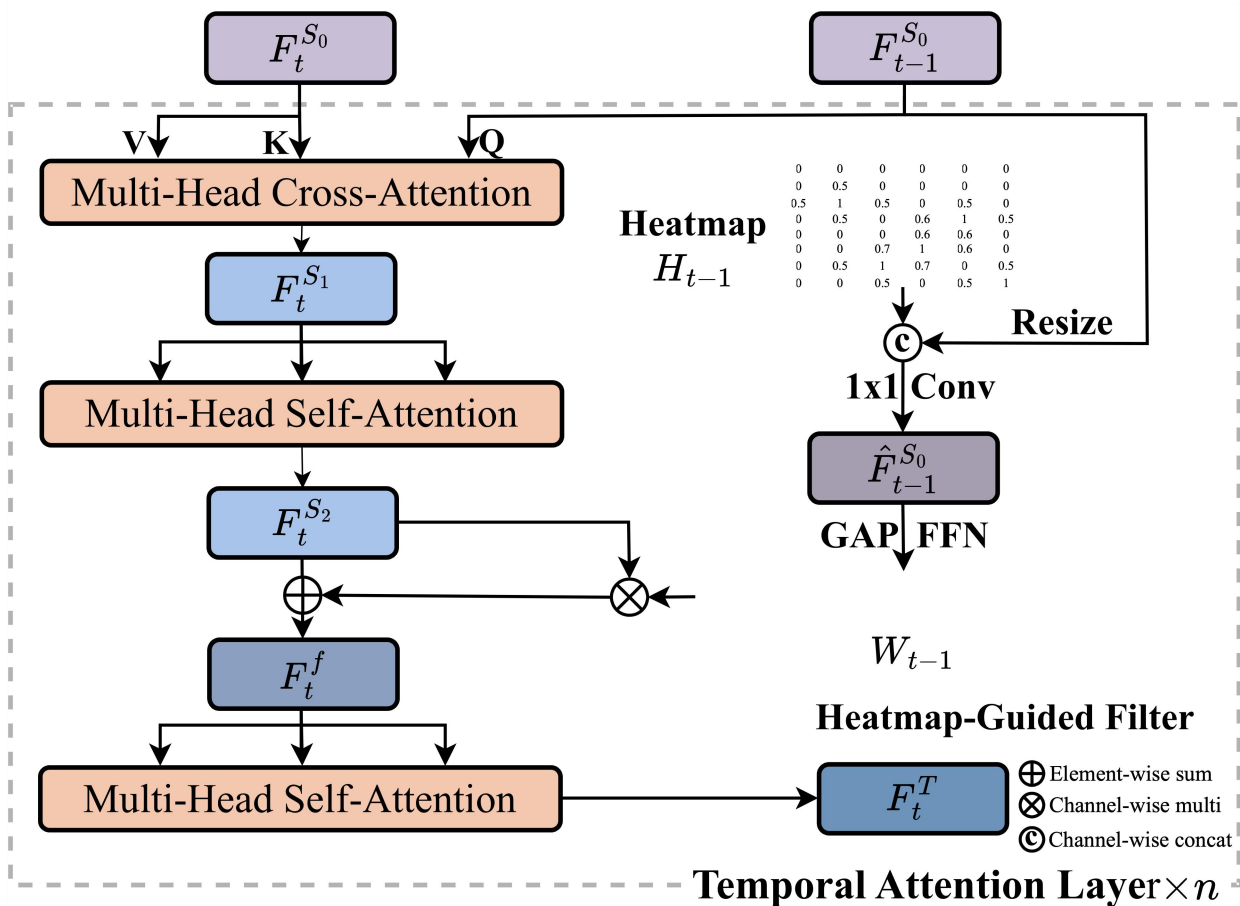
$$F_t^0 = F_t + PosEncod$$

$$F_t^{i+1} = \text{Norm} (F_t^i + \text{MultiHead}(F_t^i, F_t^i, F_t^i)), i = 0, 1,$$

$$F_t^{S_0} = \text{Norm} (F_t^2 + \text{MultiHead}(F_t^2, F_t^2, F_t^2))$$

- 增强特征向量以提升空间位置感知
- 不仅关注个体目标的空间定位，还关注目标之间的相互作用。

基于热力图的时间注意力模块



$$F_t^{S_1} = \text{Norm} (F_t^{S_0} + \text{MultiHead}(F_{t-1}^{S_0}, F_t^{S_0}, F_t^{S_0}))$$

$$F_t^{S_2} = \text{Norm} (F_t^{S_1} + \text{MultiHead}(F_t^{S_1}, F_t^{S_1}, F_t^{S_1}))$$

$$\hat{F}_{t-1}^{S_0} = \mathcal{F} (\text{Cat} (H_{t-1}, \text{Resize}(F_{t-1}^{S_0})))$$

$$W_{t-1} = \text{FFN} (\text{GAP}(\hat{F}_{t-1}^{S_0}))$$

$$F_t^f = F_t^{S_2} + F_t^{S_2} \times W_{t-1}$$

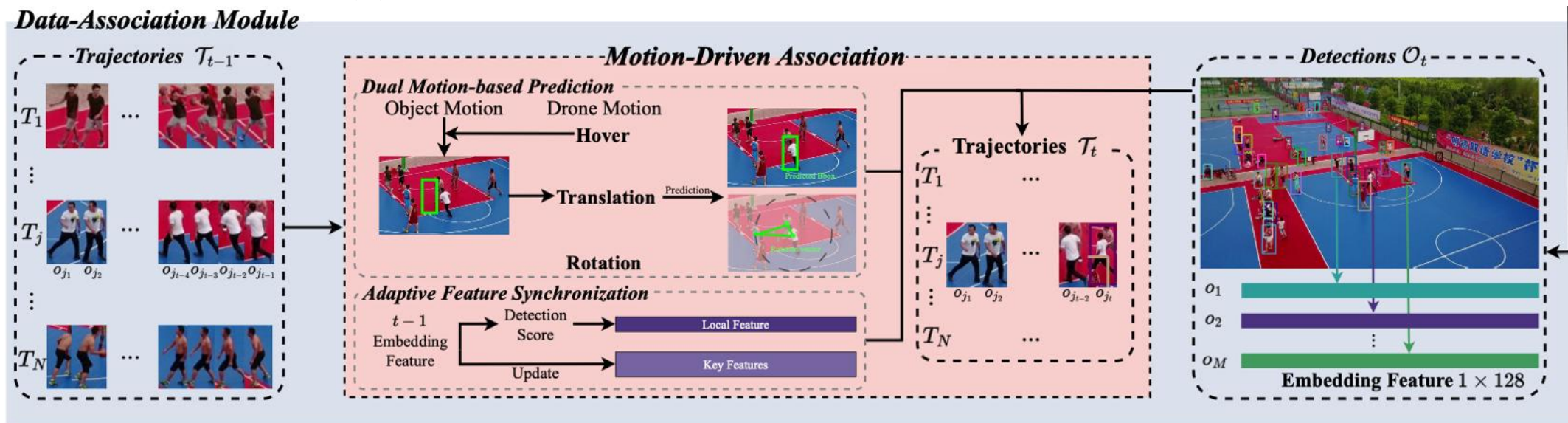
$$F_t^T = \text{Norm} (F_t^f + \text{MultiHead}(F_t^f, F_t^f, F_t^f))$$

热图引导的时间注意力层关键地处理目标特征随时间的演变，在无人机追踪中尤为重要，因为运动模糊和遮挡影响了时间上下文。

使用双域注意力与不使用双域注意力的特征图对比



基于运动的目标关联和匹配



Dual Motion-based Prediction

- Decomposing Drone Flight Patterns

- Hover

- $x_k = [x_c(k), y_c(k), w(k), h(k), \dot{x}_c(k), \dot{y}_c(k), \dot{w}(k), \dot{h}(k)]^T$

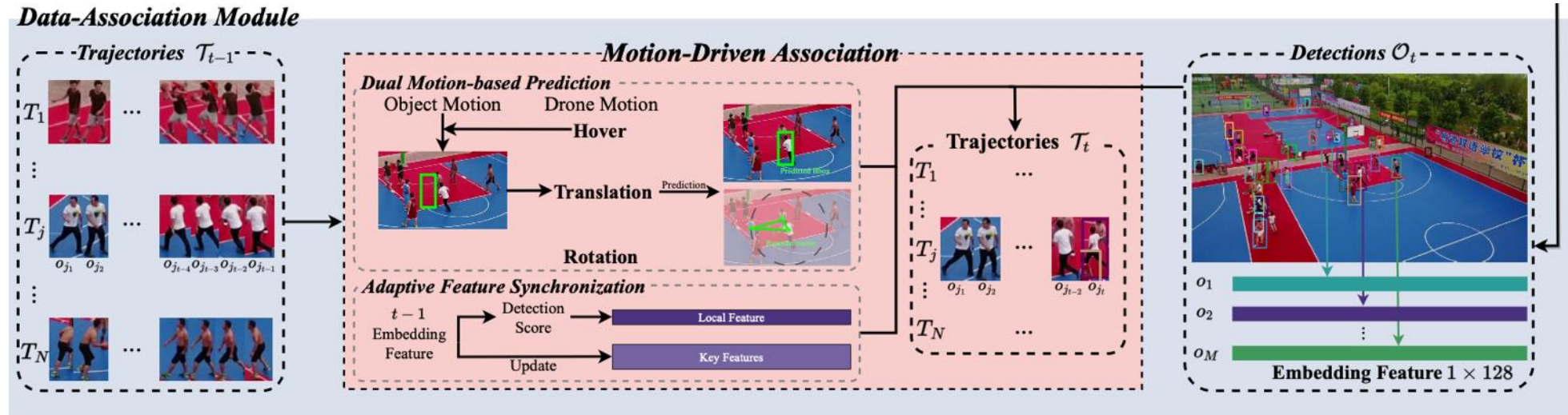
- Translation

- $A_{k-1}^k = [M_{2 \times 2} | T_{2 \times 1}] = \begin{bmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \end{bmatrix}$

- Rotation

- Key geometric features in the neighborhood graph

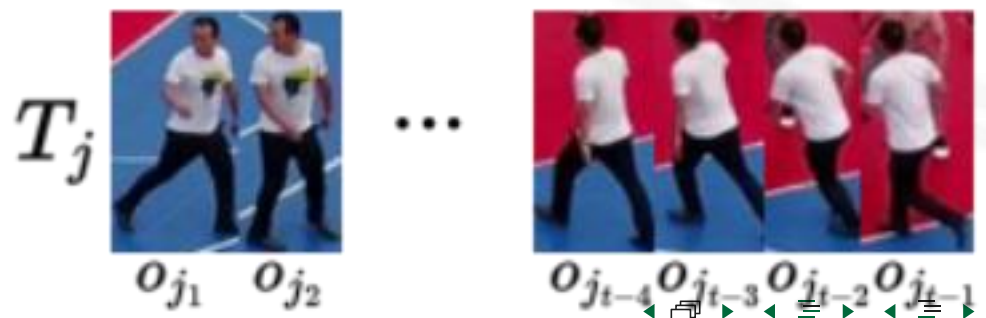
Motion-Driven Association



Adaptive Feature Synchronization

- Local Feature
 - $e_t = \alpha_t e_{t-1} + (1 - \alpha_t) e_t$
 - $\alpha_t = \alpha_f + (1 - \alpha_f) \left(1 - \frac{\det^{-\sigma}}{1 - \sigma}\right)$

- Key Feature
 - $B = B \cup \{e_t \mid D(e_t, e_{t-1}) > \Delta\}$



实验结果

在VisDrone-2019, UATDT数据集上与现有的MOT方法相比, **具有更高的IDF1, MOTA
准确性更好。**

Dataset	Method	Pub&Year	IDF1↑	MOTA↑	MOTP↑	MT↑	ML↓	FP↓	FN↓	IDs↓
VisDrone2019-MOT	SiamMOT [64]	CVPR2021	48.3	31.9	73.5	-	-	24123	142303	862
	MOTR [38]	ECCV2022	41.4	22.8	72.8	272	825	28407	147937	959
	ByteTrack [60]	ECCV2022	40.8	25.1	72.4	446	1099	34044	194984	1590
	OC-SORT [56]	CVPR2023	50.4	39.6	73.3	-	-	14631	123513	986
	STDFormer [65]	TCSVT2023	57.1	45.9	77.9	684	538	21288	101506	1440
	UAVMOT [24]	CVPR2022	51	36.1	74.2	520	574	27983	115925	2775
	FOLT [48]	MM2023	56.9	42.1	77.6	-	-	24105	107630	800
	GLOA [47]	J-STARS2023	46.2	39.1	76.1	581	824	18715	158043	4426
DroneMOT	Ours		58.6	43.7	71.4	689	397	41998	86177	1112
UATDT	DeepSORT [32]	ICIP2017	58.2	40.7	73.2	595	338	44868	155290	2061
	SiamMOT [64]	CVPR2021	61.4	39.4	76.2	-	-	46903	176164	190
	ByteTrack [60]	ECCV2022	59.1	41.6	79.2	-	-	28819	189197	296
	OC-SORT [56]	CVPR2023	64.9	47.5	74.8	-	-	47681	148378	288
	UAVMOT [24]	CVPR2022	67.3	46.4	72.7	624	221	66352	115940	456
	FOLT [48]	MM2023	68.3	48.5	80.1	-	-	36429	155696	338
	GLOA [47]	J-STARS2023	68.9	49.6	79.8	626	220	55822	115567	433
	DroneMOT	Ours		69.6	50.1	74.5	638	178	57411	112548

消融实验结果：所提的DIA模块和MDA模块都有明显效果

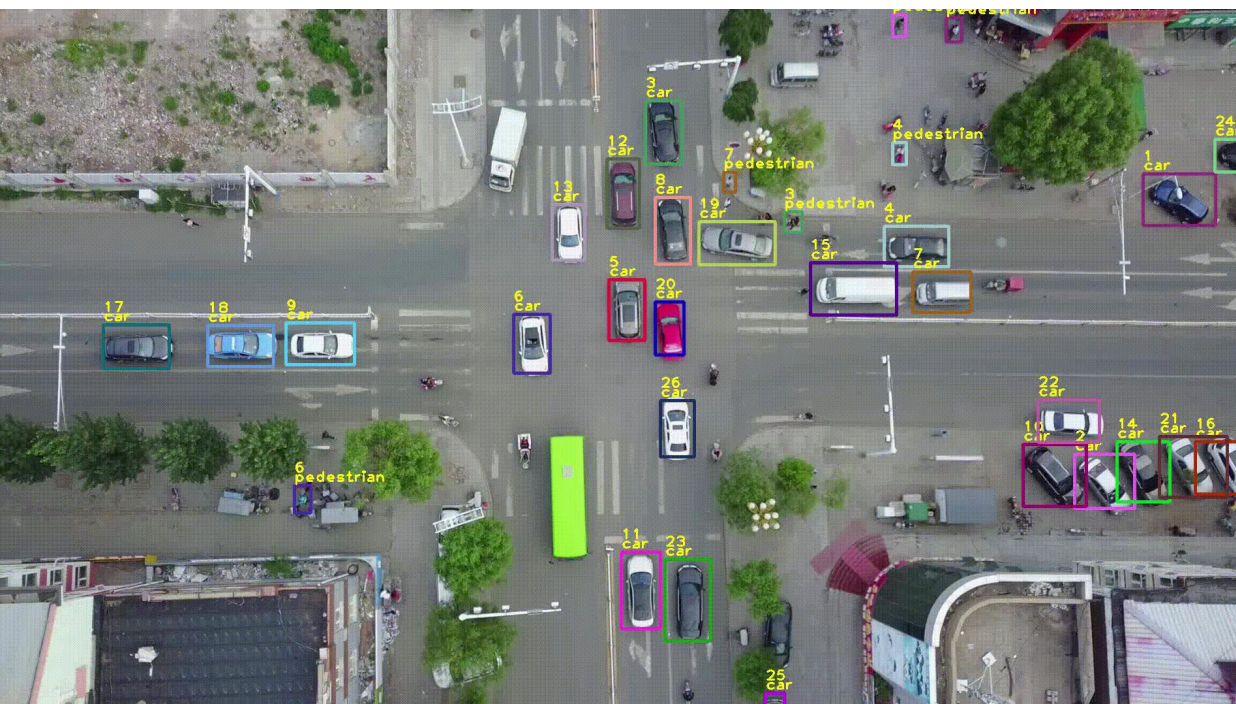
TABLE II
ABALATION STUDY ON VISDRONE2019-MOT VALIDATION SET.

Baseline	DIA	MDA	MOTA(%)	IDs	IDF1(%)
✓			29.7	1509	38.3
✓	✓		33.4	1407	45.1
✓		✓	32.4	406	48.9
✓	✓	✓	34.3	218	53.4

TABLE III
ANALYSIS OF THE EFFECTIVENESS OF MDA MODULE. THE BASELINE USES THE KALMAN FILTER AND EMA TO UPDATE THE FEATURE.

Motion model	Appearance model	IDs	IDF1	IDP	IDR
-	-	1407	45.1	48.6	42.1
DMP	-	229	52.8	57.8	48.6
-	AFS	690	46.5	52.8	41.5
DMP	AFS	218	53.4	43.0	52.8

可视化对比: DroneMOT and UAVMOT



DroneMOT 我们的方法

无人机旋转 情况下跟踪的更为准确



UAVMOT

无人机旋转

DroneMOT夜间多车追踪效果



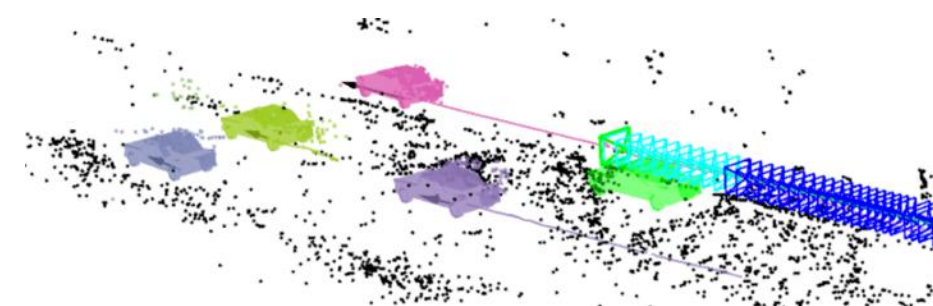
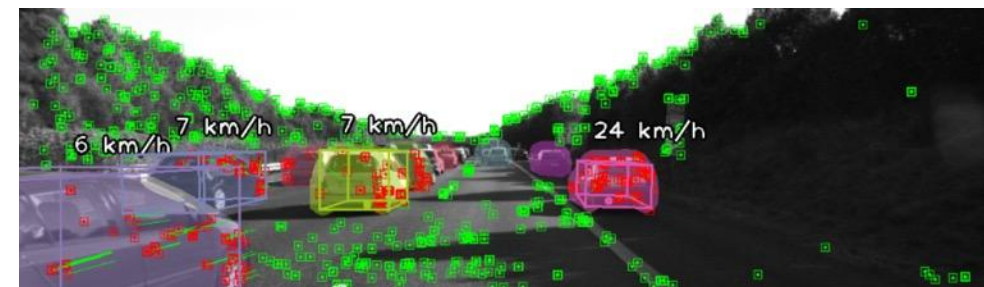
GSLAMOT: A Tracklet and Query Graph-based Simultaneous Locating, Mapping, and Multiple Object Tracking System

GSLAMOT: 基于轨迹图与查询图的同步定位 建图与多目标追踪系统

投稿于ACM MM2024, 人工智能领域顶会
王硕、王永才等
中国人民大学信息学院

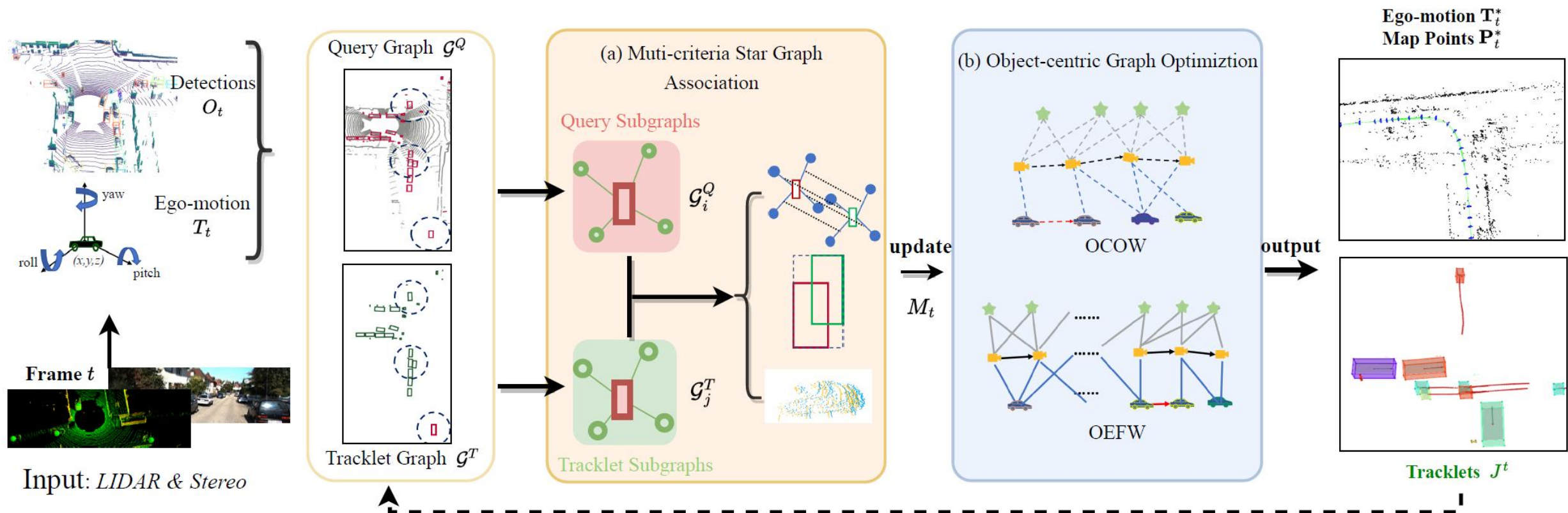
2. GSLAMOT: 实现无人车的同步定位建图与多目标追踪系统

- 系统输入
 - 无人车连续采集雷达与双目视觉相机数据
- 系统输出
 - 无人车自身的位姿轨迹
 - 建立静态语义环境地图
 - 检测并追踪周边的移动车辆



Locating + Mapping + Tracking

系统架构：同步定位建图与多目标追踪系统



输入

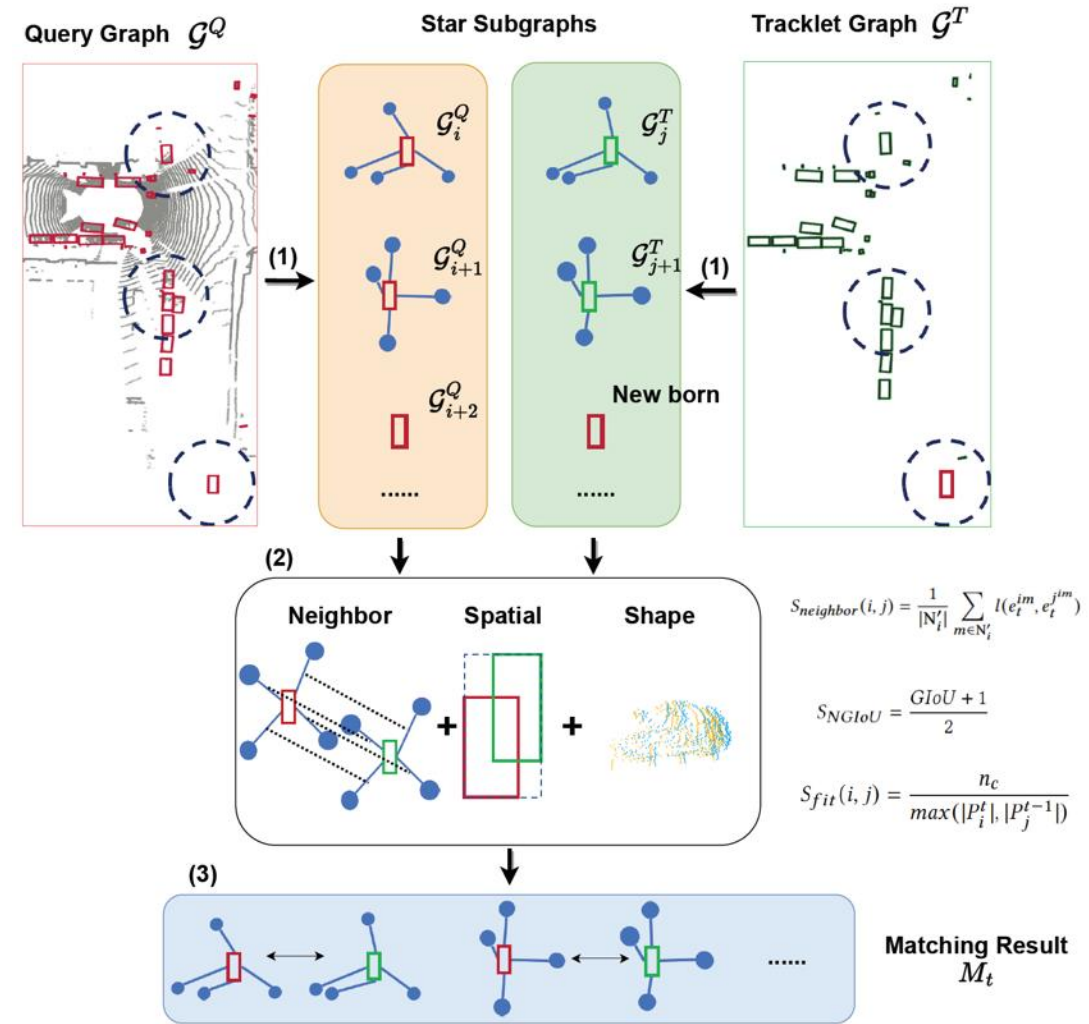
目标检测

当前帧查询历史目标轨迹图

自身轨迹、建图、
多目标跟踪结果

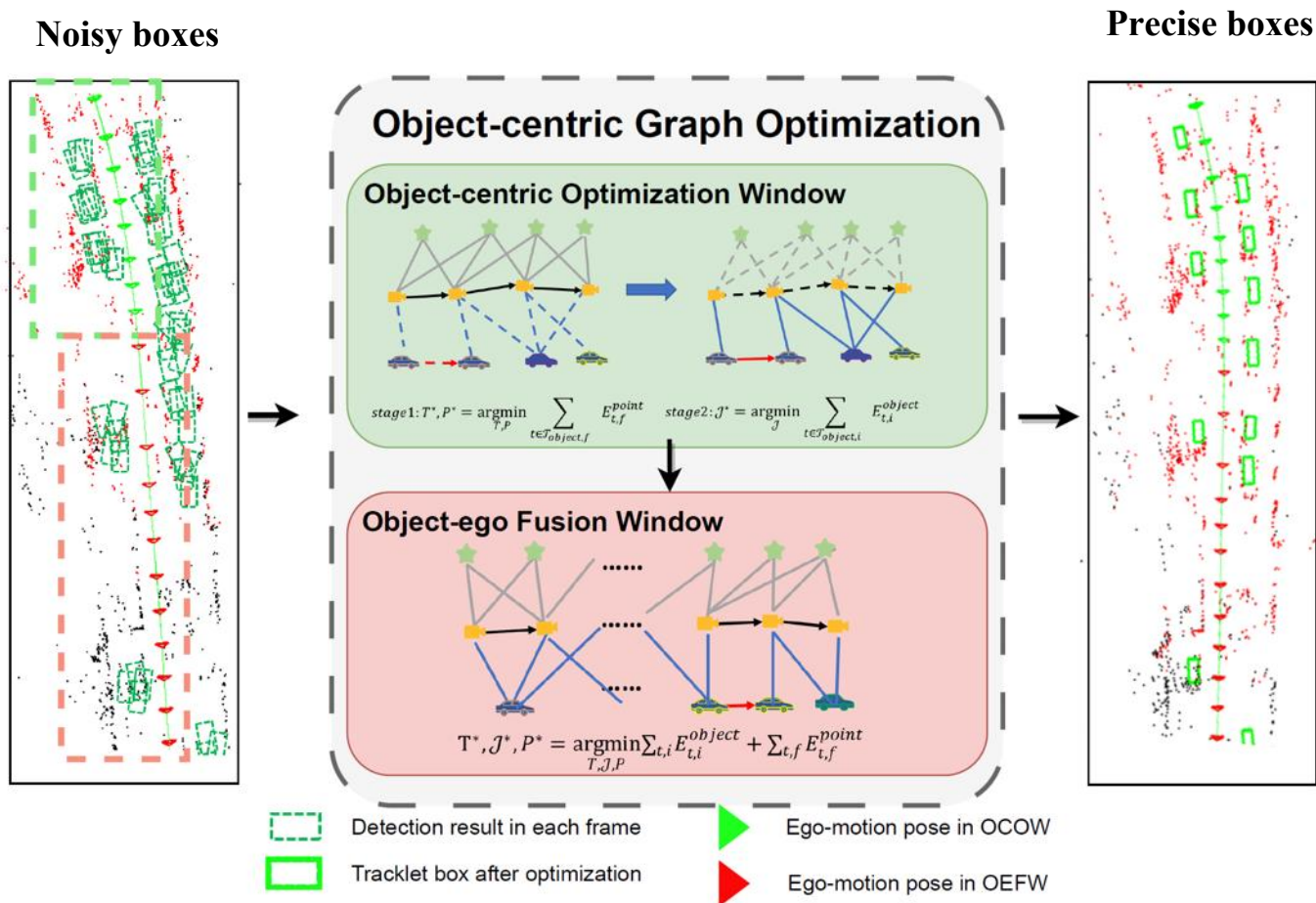
多种相似性的查询图(Query Graph)与历史轨迹图(Tracklet Graph)匹配

- 我们为当前帧的检测创建**查询图(Query Graph)**，并为地图中的轨迹创建**轨迹图(Tracklet Graph)**。
- 每个检测和轨迹都分别被分配一个**星形子图**。
- 我们通过评估它们的**星形子图的邻居、空间和形状一致性**来**匹配检测和轨迹**。



以目标为中心的图优化方法

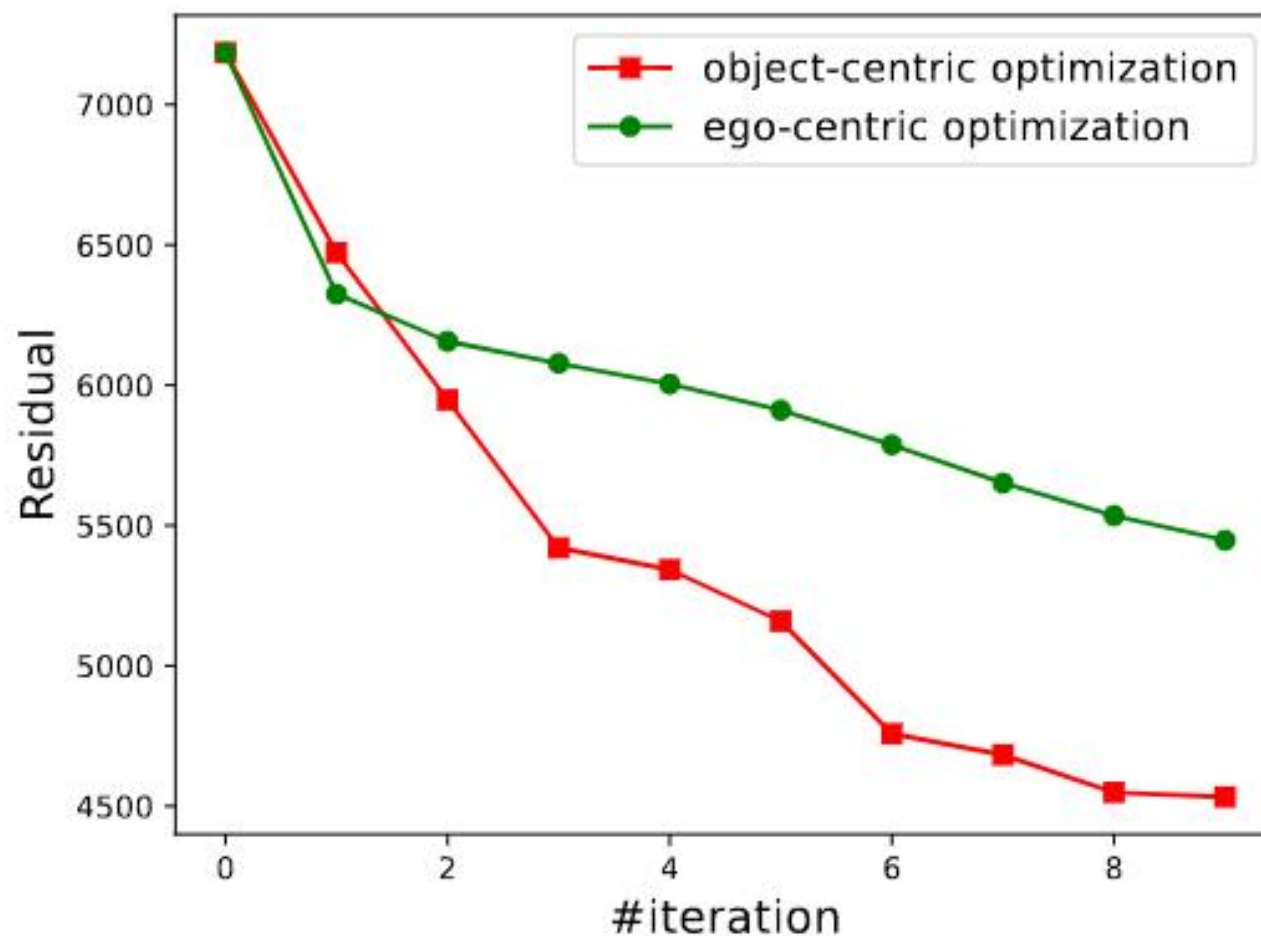
- 以目标为中心的图优化窗口 (OCOW)
 - 阶段1: 利用静态环境地标估计车辆的自身运动。
 - 阶段2: 固定车辆自身运动, 来优化移动目标位姿。
- 目标-自我融合窗口 (OEFW):
 - 一个紧密耦合的优化策略, 联合优化自我运动姿势、地图点和轨迹姿势。



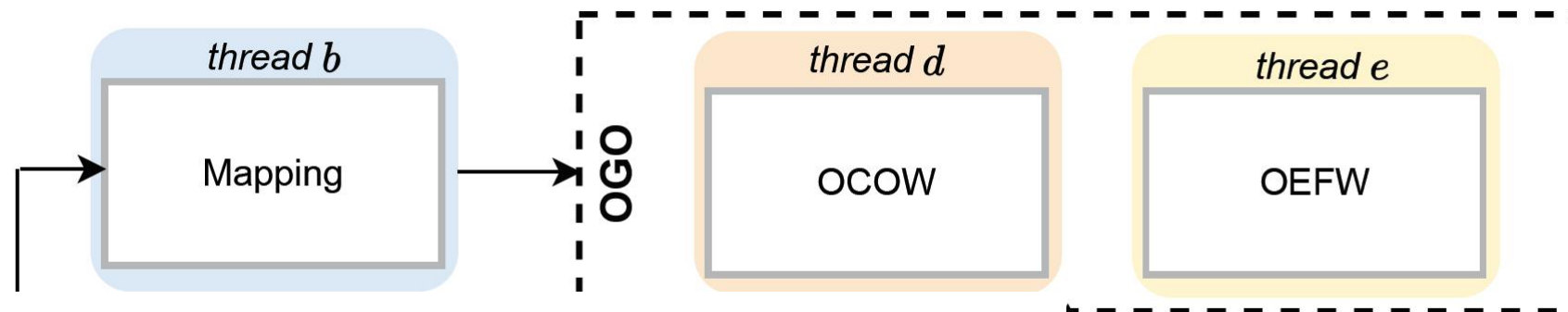
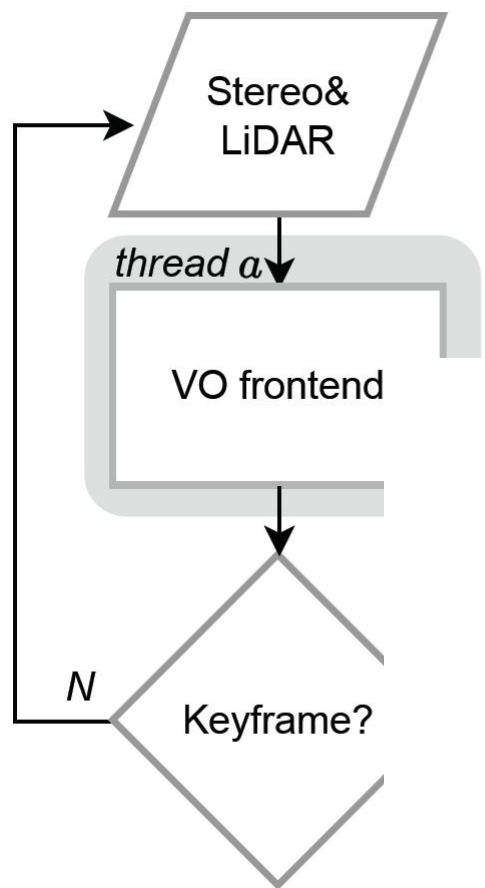
以目标为中心的图优化方法

我们提的以目标为中心的优化和经典的优化方法的收敛残差对比。

误差更小、收敛更快。

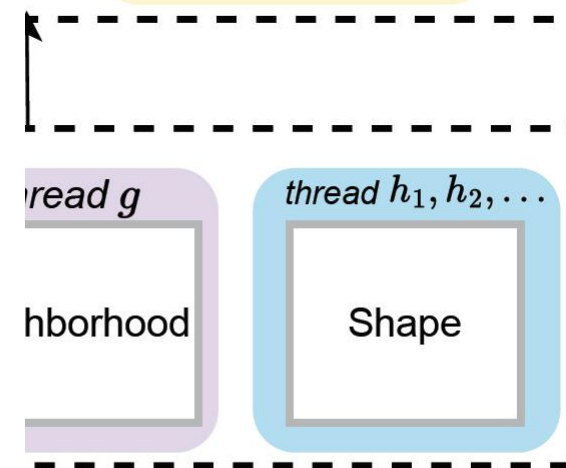


多线程的系统实现



Module	Time Per Frame(s)↓
Ego-motion Tracking	0.031
Mapping	0.027
3D Detection†	0.022
MSSA†	0.024
OGO†	0.037
Total	0.082

†: only for keyframes.



实验数据集



我们自己构建的高密度交通流数据集

实验结果：定位与多目标追踪效果

在KITTI数据集上结果显著好于所有现有方法

KITTI Seq.	00		01		02		03		04		05		06		07		08		Average	
Metrics(m)	RPE	APE	RPE	APE	RPE	APE	RPE	APE	RPE	APE	RPE	APE	RPE	APE	RPE	APE	RPE	APE	RPE	APE
ORB_SLAM3[5]	2.09	1.46	7.52	<u>12.70</u>	2.3	3.5	0.84	1.44	<u>0.6</u>	<u>0.25</u>	0.91	0.93	0.92	0.99	0.49	0.49	<u>3.06</u>	3.06	2.08	<u>2.76</u>
DSP-SLAM[33]	1.09	<u>1.10</u>	3.87	12.06	<u>0.94</u>	0.89	1.28	0.47	0.64	0.73	0.53	<u>0.46</u>	0.81	0.42	0.5	0.48	3.17	11.99	1.40	3.18
DynaSLAM[3]	1.05	1.28	<u>3.75</u>	21.13	1.1	<u>0.91</u>	0.68	1.43	0.73	0.82	0.64	1.52	0.8	1.35	0.51	0.78	<u>3.06</u>	10.41	<u>1.36</u>	4.40
VDO-SLAM[42]	<u>1.02</u>	1.44	3.80	13.79	0.98	0.99	<u>0.79</u>	0.83	0.61	<u>0.25</u>	<u>0.59</u>	0.49	<u>0.75</u>	0.63	0.49	0.52	3.34	9.76	1.50	3.19
GSLAMOT(Ours)	1.01	1.02	3.69	13.1	0.92	<u>0.91</u>	0.68	<u>0.57</u>	0.56	0.23	0.53	0.41	0.70	<u>0.44</u>	0.48	0.43	3.05	<u>3.16</u>	1.29	2.25

在Waymo数据集上3D目标检测效果优于现有方法.

Method	MOTA(L1)↑	MOTA(L2)↑	Mismatch↓	MOTA(L2)↑		
				vehicle	pedestrian	cyclist
AB3DMOT[34]	-	-	-	40.1	33.7	50.39
ProbTrack[7]	48.26	45.25	1.05	54.06	48.10	22.98
CenterPoint[37]	58.35	55.81	0.74	59.38	56.64	60.0
SimpleTrack[25]	59.44	56.92	0.36	56.12	57.76	56.88
BOTT[46]	59.67	<u>57.14</u>	<u>0.35</u>	59.49	58.82	60.41
TrajectoryF[6]	-	-	-	59.7	61.0	60.6
GSLAMOT	<u>59.69</u>	57.10	0.33	<u>60.45</u>	60.02	60.33
GSLAMOT*	59.75	57.20	0.33	60.47	<u>60.23</u>	<u>60.45</u>

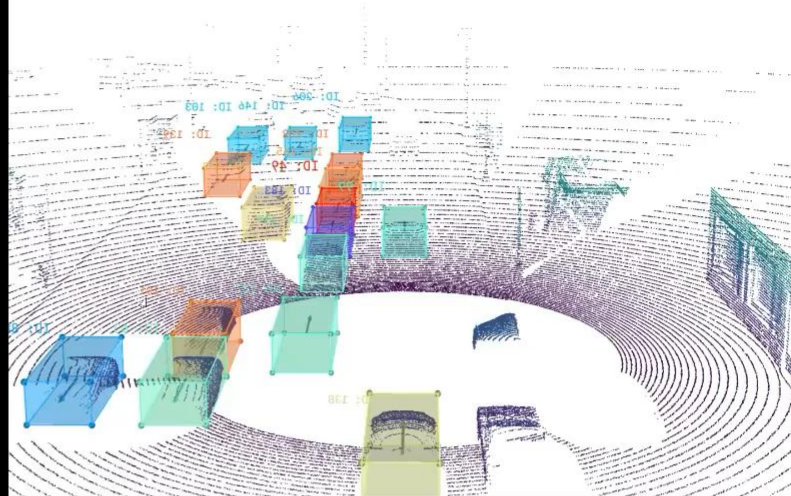
GSLAMOT: The ego-motion poses are estimated by odometry front-end.

GSLAMOT*: The groundtruth ego-motion poses are given as other MOT algorithms.

录像展示

Speed X 10

多目标追踪效果，
不同颜色代表追踪
的不同目标。



- Green box: tracklet
- Red point: local map
- Black point: map point

SAHNet :考虑语义信息与智能体差异的协同3D目标检测

RoCo: 基于迭代目标匹配与位姿矫正的鲁棒协同感知

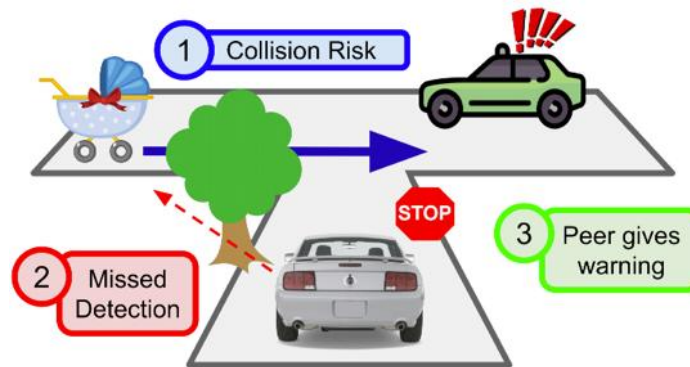
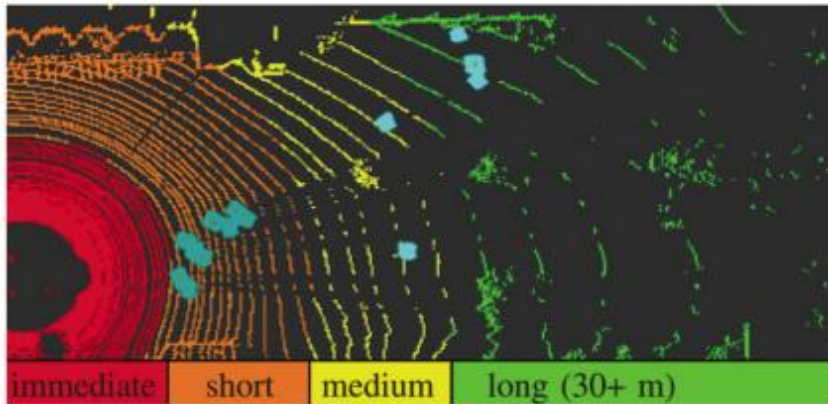
投稿于TIV2024, ACM MM2024

黄哲、王永才等

中国人民大学信息学院

Single-Agent Perception

- Single-vehicle self-driving has been intensively studied^[1]
- Long-range perception is challenging due to the sparse measurements of 3D sensor^[2]
- Individual viewpoint suffers from frequent occlusions^[2]

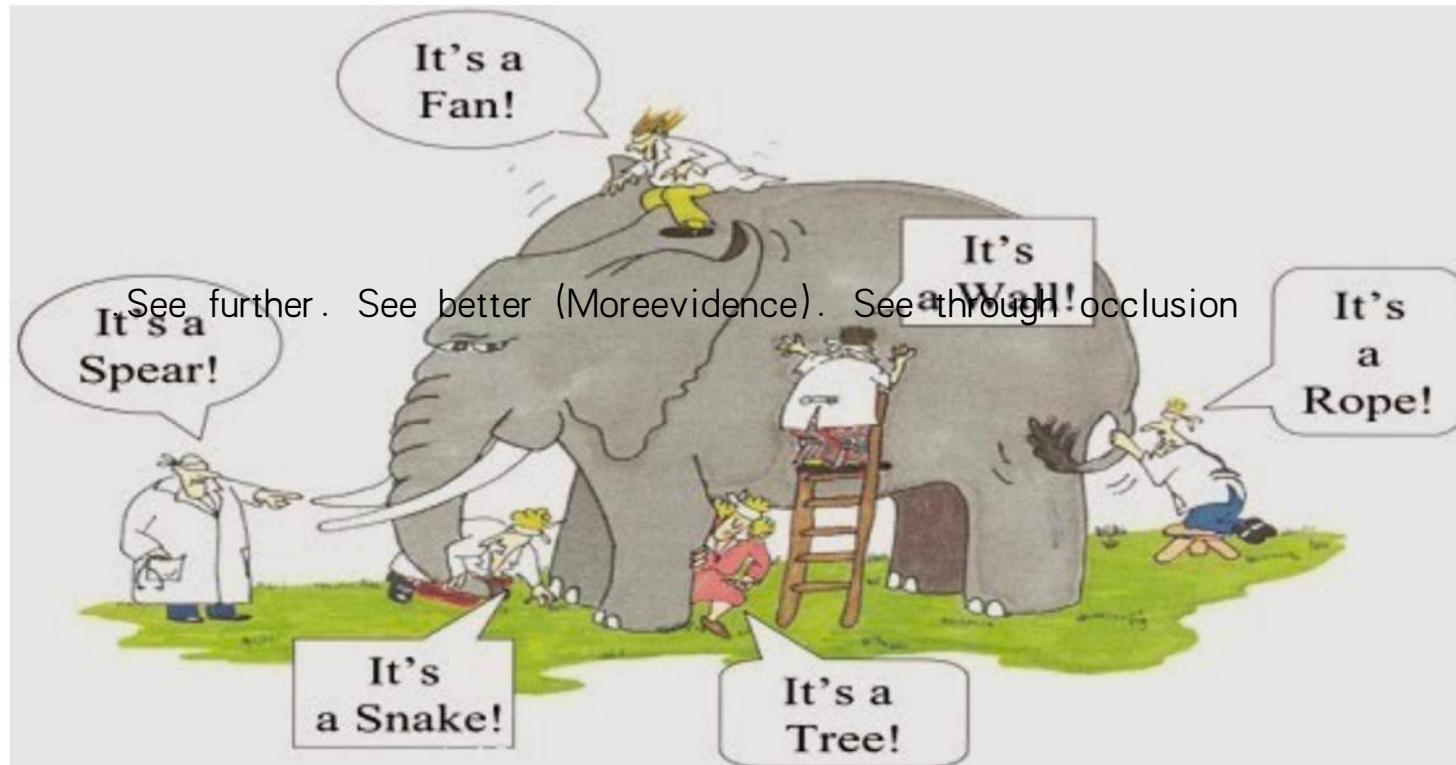


[1] Li Y, Ibanez-Guzman J. Lidar for autonomous driving: The principles, challenges, and trends for automotive lidar and perception systems[J]. IEEE Signal Processing Magazine, 2020, 37(4): 50-61.

[2] Wang T H, Manivasagam S, Liang M, et al. V2vnet: Vehicle-to-vehicle communication for joint perception and prediction[C]//Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part II 16. Springer International Publishing, 2020: 605-621.

Collaborative Perception

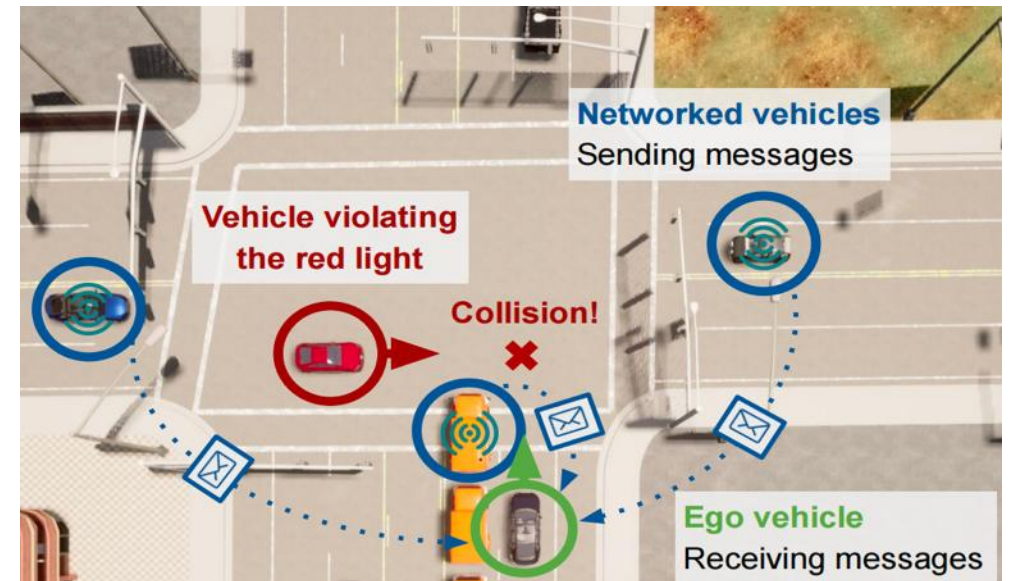
What you see is what you get ?



Collaboration! Holistic view!

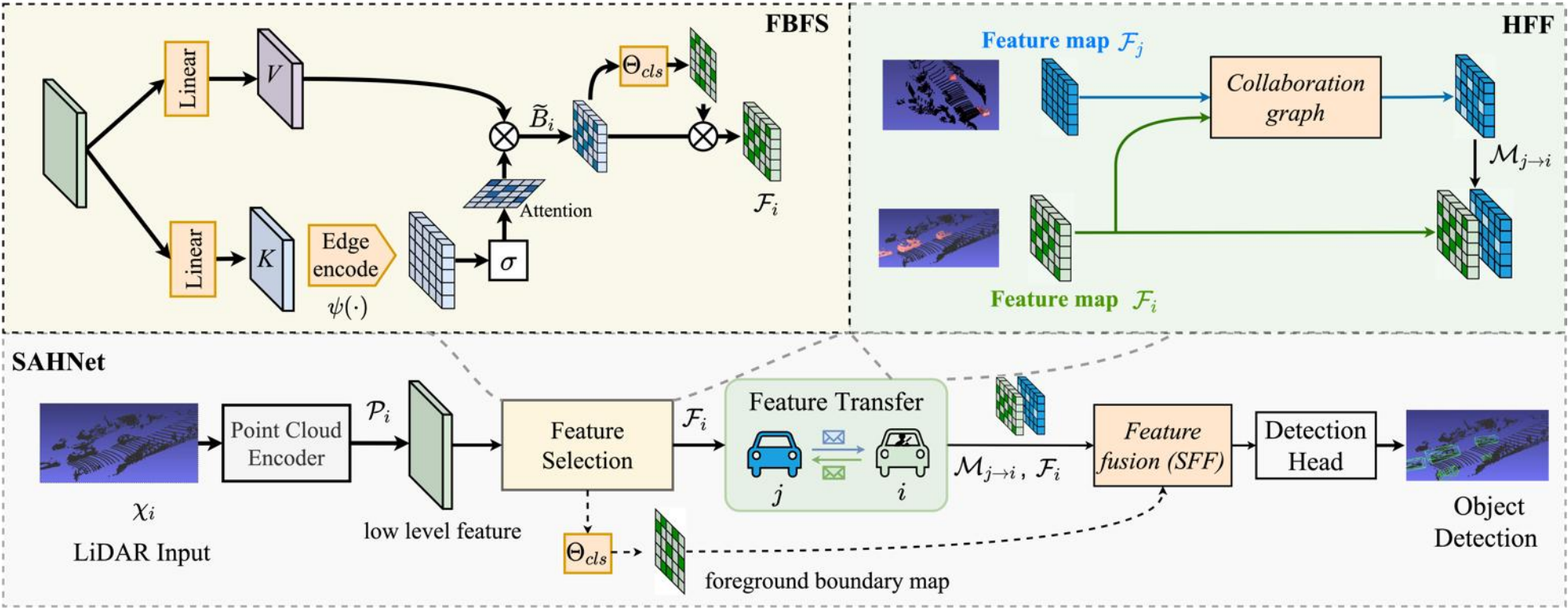
Collaborative Perception

- ✓ See further
- ✓ See better (More evidence).
- ✓ See through occlusion



SAHNet: Collaborative Object Detection Method

Core idea: Exploring spatial heterogeneity of perceptual information
 Messages should be spatially sparse, yet perceptually critical



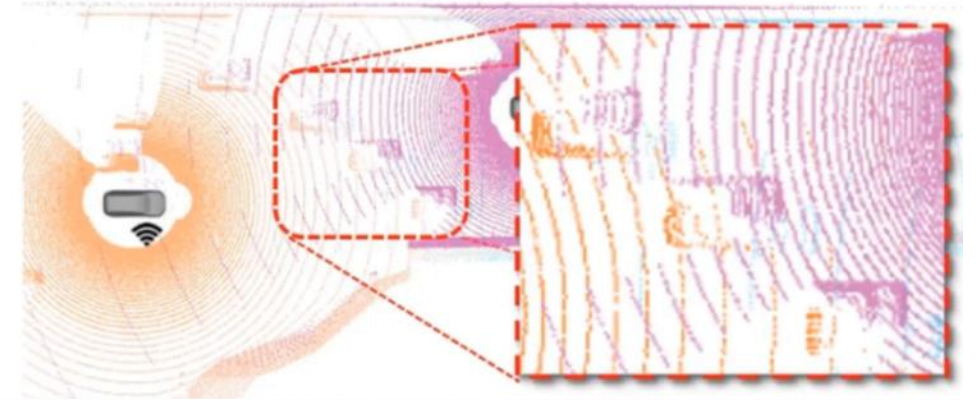
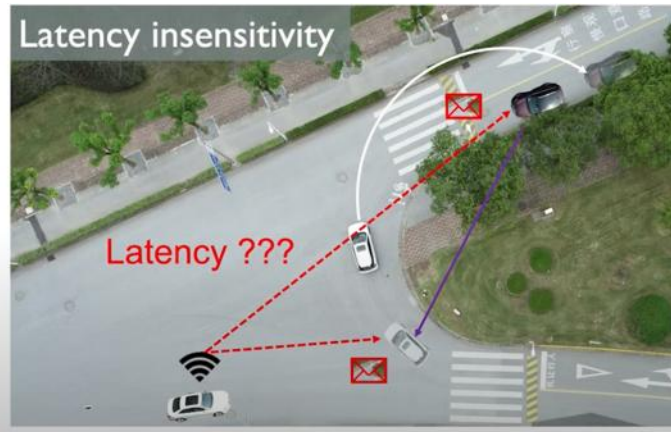
Current work: Pose error and time delay in communication

Real-world communication among agents :

- Congestion
- Heavy computation
- Interruptions
- Lack of calibration



delays and misalignments

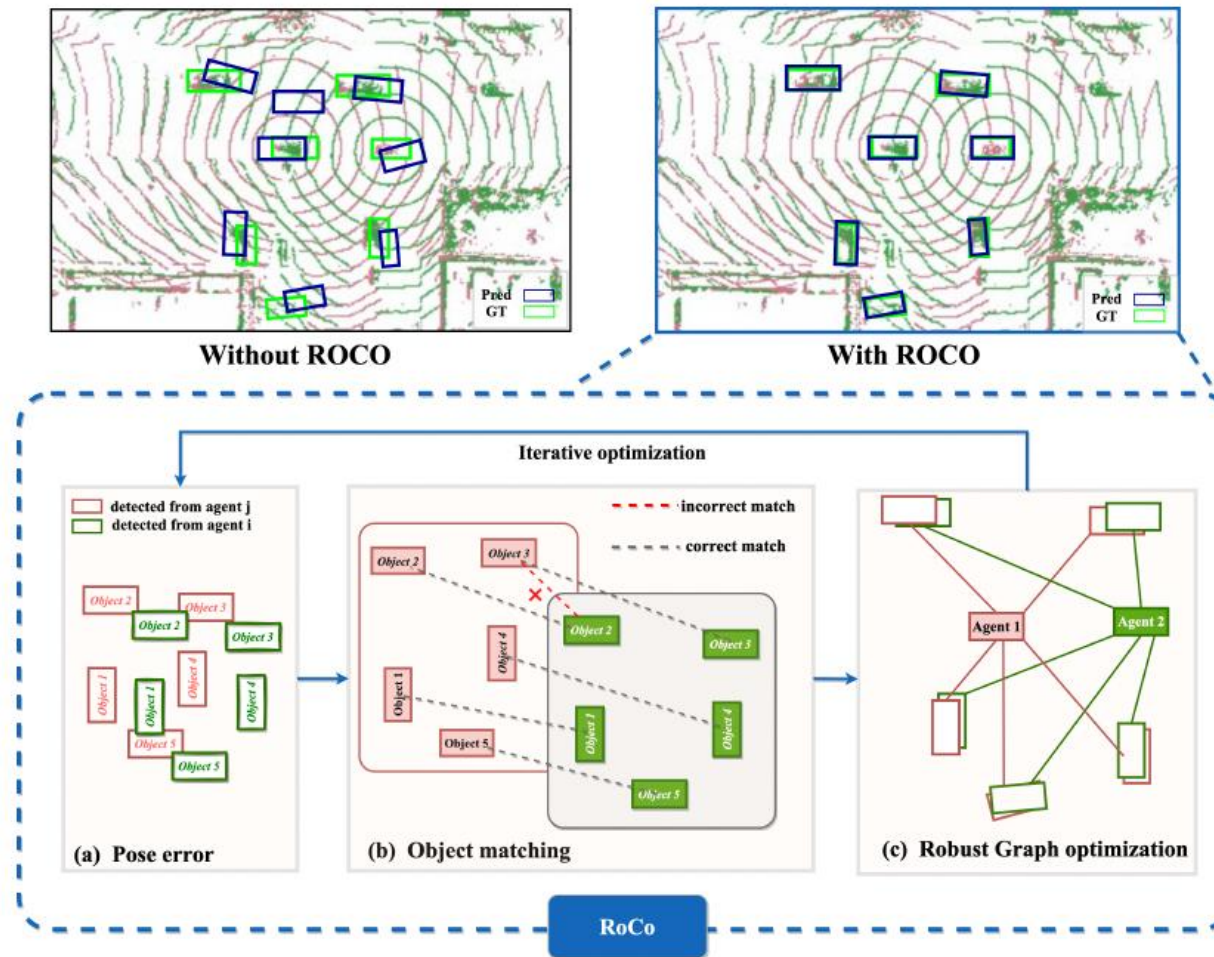


(a) A collision caused by latency

(b) Pose errors degrade collaboration performance

Collaborative perception--RoCo

To minimize the impact of pose errors between agents, RoCo was proposed





*RoCo: Robust Collaborative Perception By Iterative Object
Matching and Pose Adjustment*

MobiSketch: 基于手机的点线面特征融合的3D语义建图

投稿于TIV2024, ACM MM2024

黄哲、王永才等

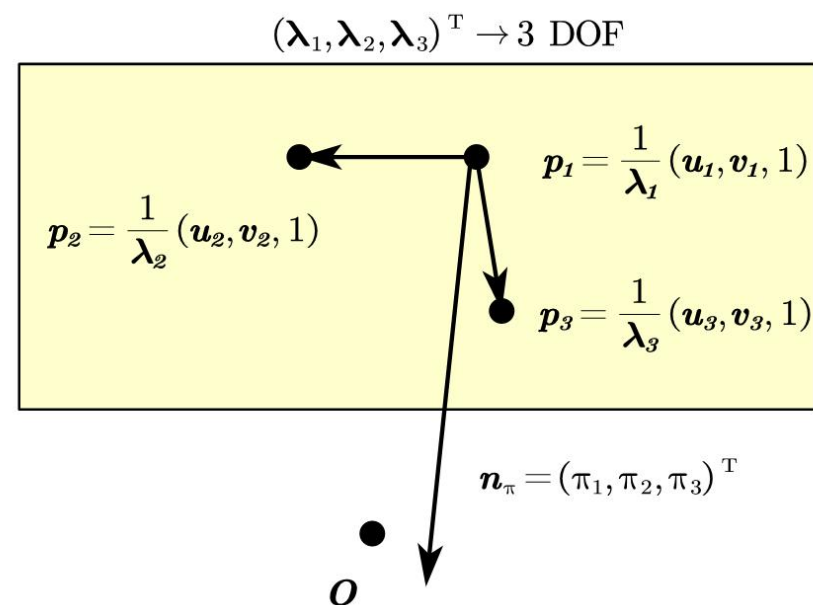
中国人民大学信息学院

面特征的表达

逆深度法-后端优化

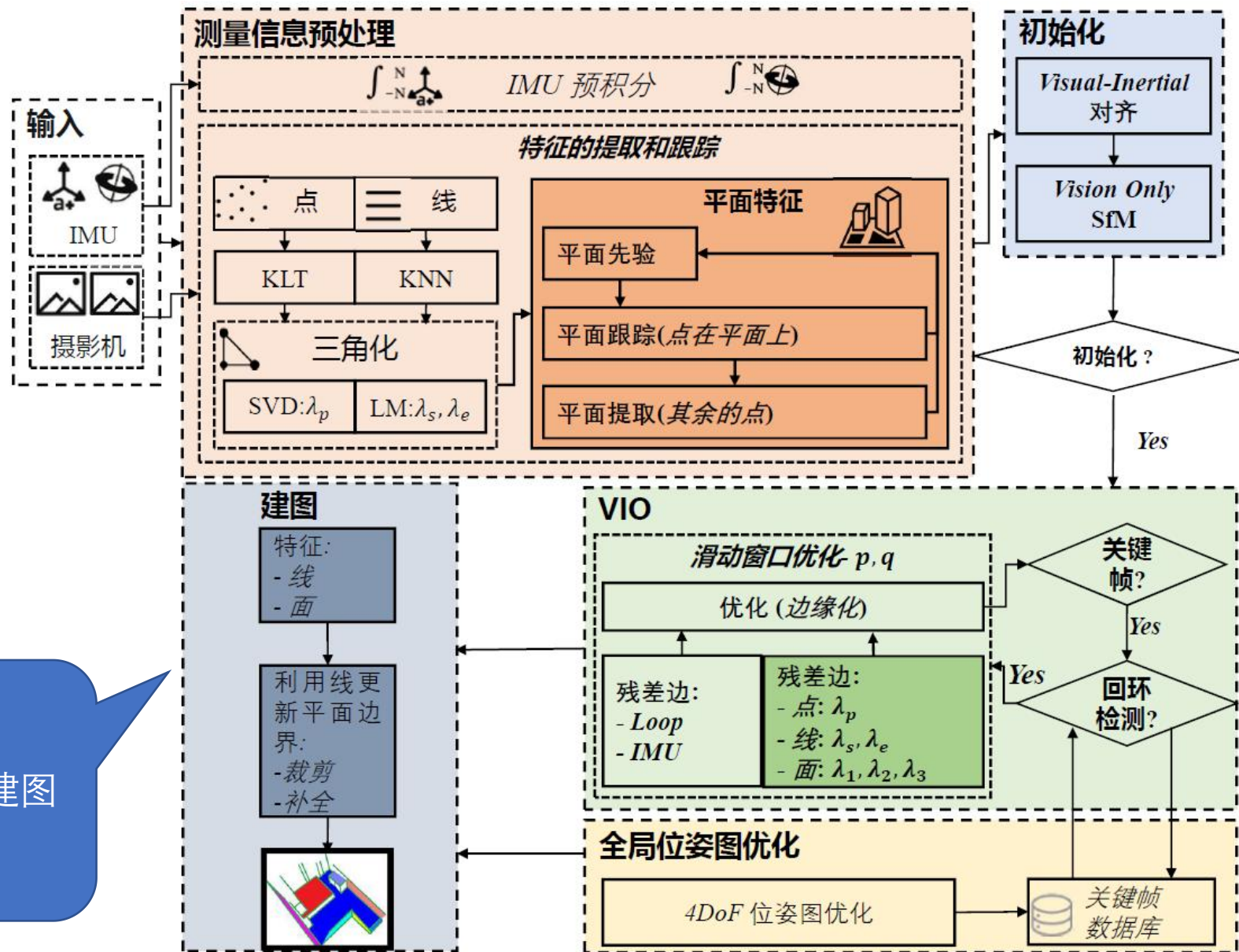
在空间中，三个非共线的点确定一个平面，这三个点在3D空间中确实具有9-DoF，这是因为每个点都可以在三个轴上独立地取值。而这个平面又对这三个点**施加点在平面上约束，最后退化为3-DoF。**

$$\mathbf{n}_\pi = \left(\frac{1}{\lambda_2} \begin{bmatrix} u_2 \\ v_2 \\ 1 \end{bmatrix} - \frac{1}{\lambda_1} \begin{bmatrix} u_1 \\ v_1 \\ 1 \end{bmatrix} \right) \times \left(\frac{1}{\lambda_3} \begin{bmatrix} u_3 \\ v_3 \\ 1 \end{bmatrix} - \frac{1}{\lambda_1} \begin{bmatrix} u_1 \\ v_1 \\ 1 \end{bmatrix} \right)$$
$$\pi_4 = -\mathbf{n}_\pi \cdot \frac{1}{\lambda_1} \begin{bmatrix} u_1 \\ v_1 \\ 1 \end{bmatrix}$$



这统一了点在1-DoF、线在2-DoF和平面在3-DoF上的逆深度表示，并且这些表示是有几何意义的。

MobiSketch算法流程



3D语义建图

定位、SLAM

实验结果


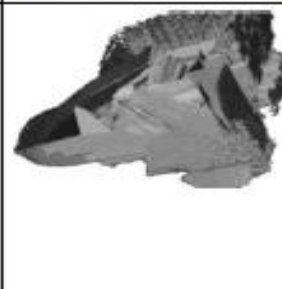
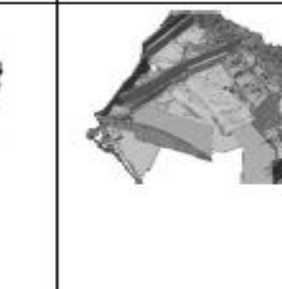
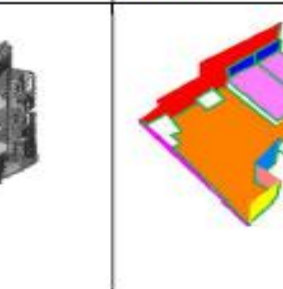

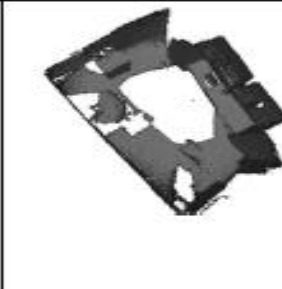
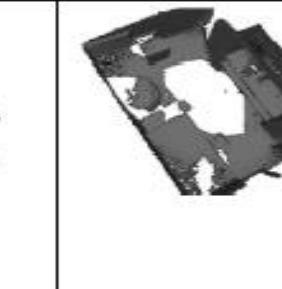
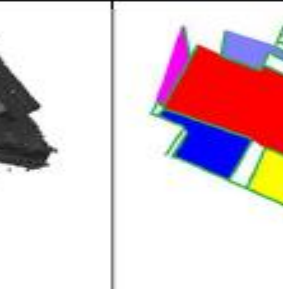
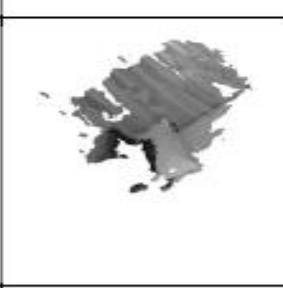

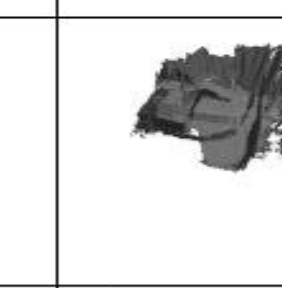
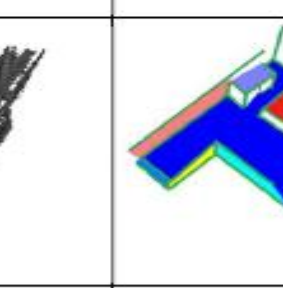


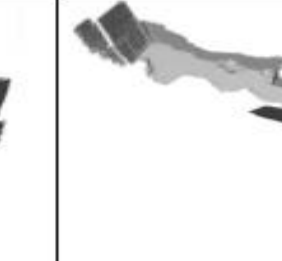
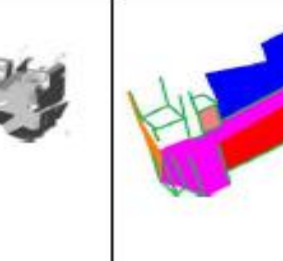
方法	VINS-MONO (点)	PL-VINS (点+线)	MobiSketch (点+线)	UV-SLAM (点+线+消失点)	PVIO (点+面)	MobiSketch (点+面)	MobiSketch (点+线+面)
MH-01-Easy	0.075	0.083	0.071	0.079	0.130	0.081	0.070
MH-02-Easy	0.070	0.072	0.070	0.052	0.285	0.075	0.054
MH-03-Medium	0.104	0.099	0.071	0.078	0.171	0.103	0.069
MH-04-Difficult	0.220	0.202	0.170	0.179	0.303	0.215	0.167
MH-05-Difficult	0.240	0.226	0.137	0.146	0.202	0.187	0.135
V1-01-Easy	0.046	0.046	0.046	0.042	0.084	0.046	0.042
V1-02-Medium	0.091	0.079	0.073	0.079	0.105	0.088	0.070
V1-03-Difficult	0.225	0.180	0.141	0.175	0.184	0.186	0.138
V2-01-Easy	0.052	0.058	0.072	0.059	0.052	0.052	0.058
V2-02-Medium	0.139	0.133	0.074	0.115	0.202	0.135	0.069
V2-03-Difficult	0.215	0.196	0.143	0.160	0.275	0.198	0.137

- (1) 添加线特征的SLAM方法的定位性能通常比添加平面特征的方法获得更好的结果。
- (2) 对线和平面采用最小逆深度表示法确实可以帮助避免不必要的误差

补充平面特征在处理线特征退化时，有助于改善定位结果。

与现有方法对比

- 同现有3D建图方法相比，MobiSketch建立起更为清晰、包含语义信息的、包含线面信息的3D空间结构地图。

方法	ColMap	ElasticFusion	InfiniTAM	MobiSketch
Dorm1				
Dorm2				
Hotel				
Living room				

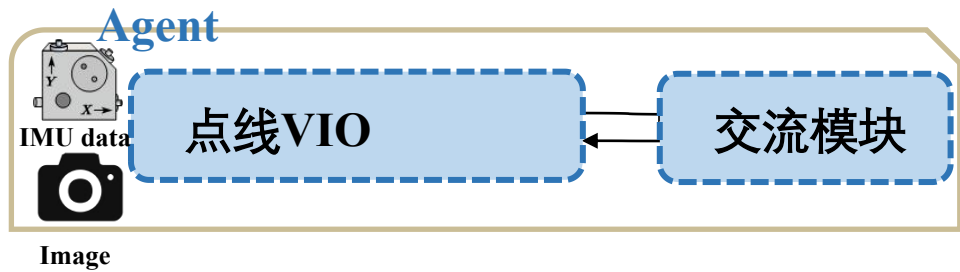
MobiSketch: Locating and Sketching on the Move Using Minimal Representation of Point, Line and Planes

CoSLAM: A Versatile Collaborative SLAM System for Mobile Phones Using Point-Line Features and Map Caching

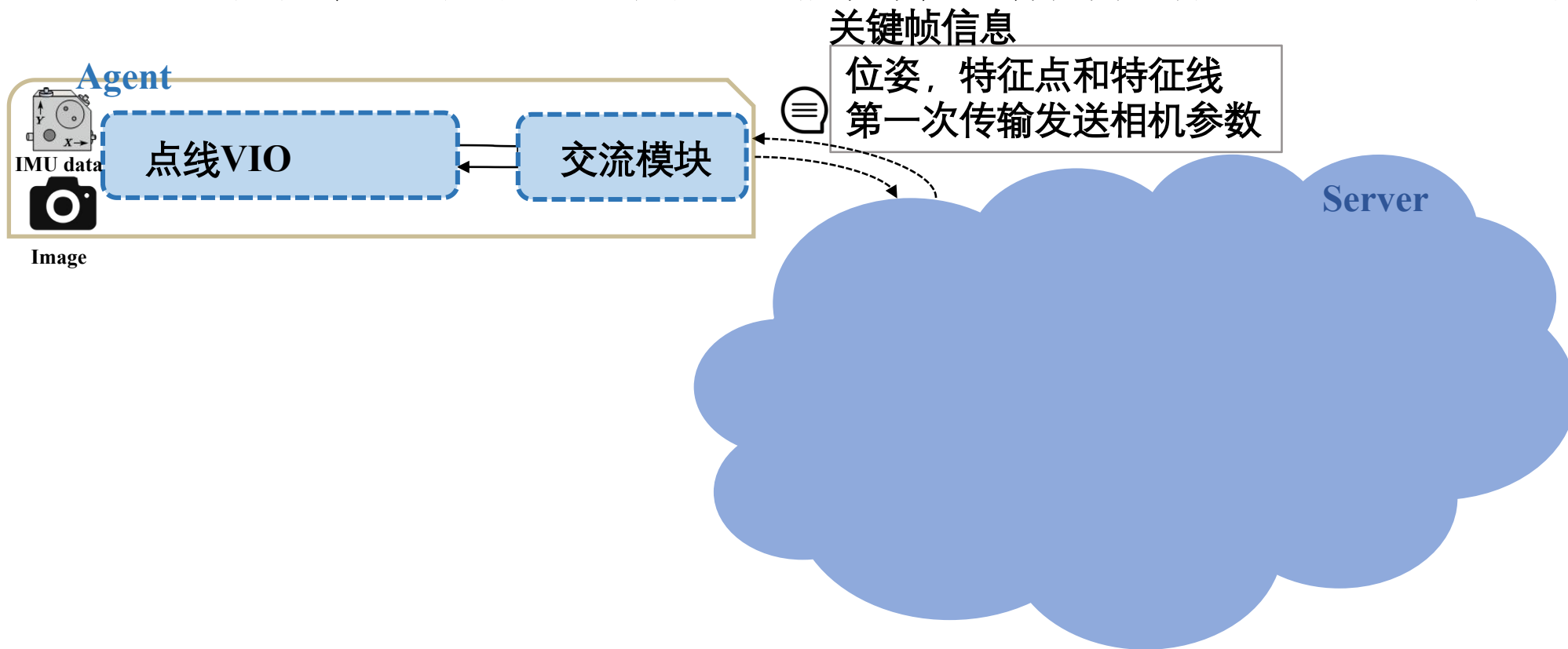
CoSLAM: 基于点线特征和缓存地图的手机协同SLAM

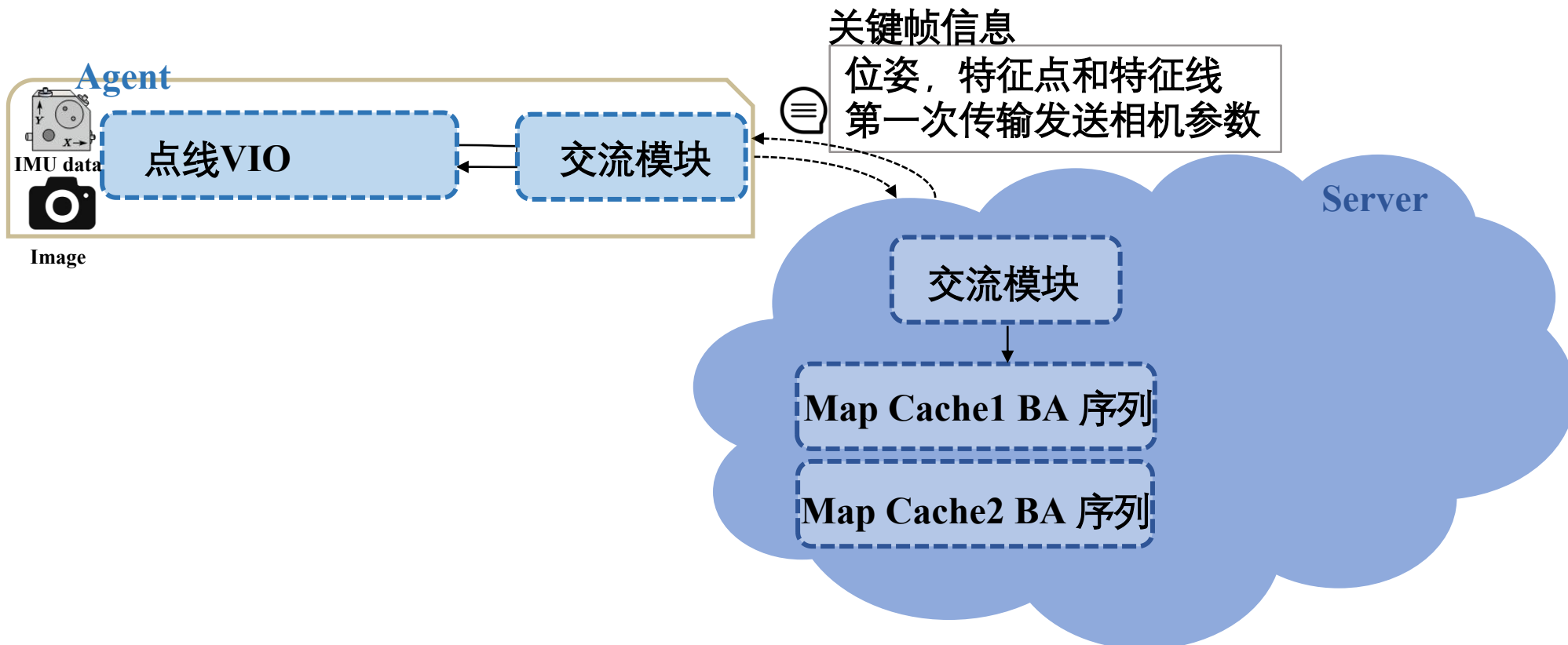
发表于 ACM MM2023, CCF A
李婉婷、王永才等
中国人民大学信息学院

CoISLAM 包括一个两阶段异步优化方法，用于缓存子图（CSG）的优化和全局地图的优化。该方法确保对代理的实时响应，同时考虑了准确性和通信效率，显著提高了协同SLAM的可扩展性

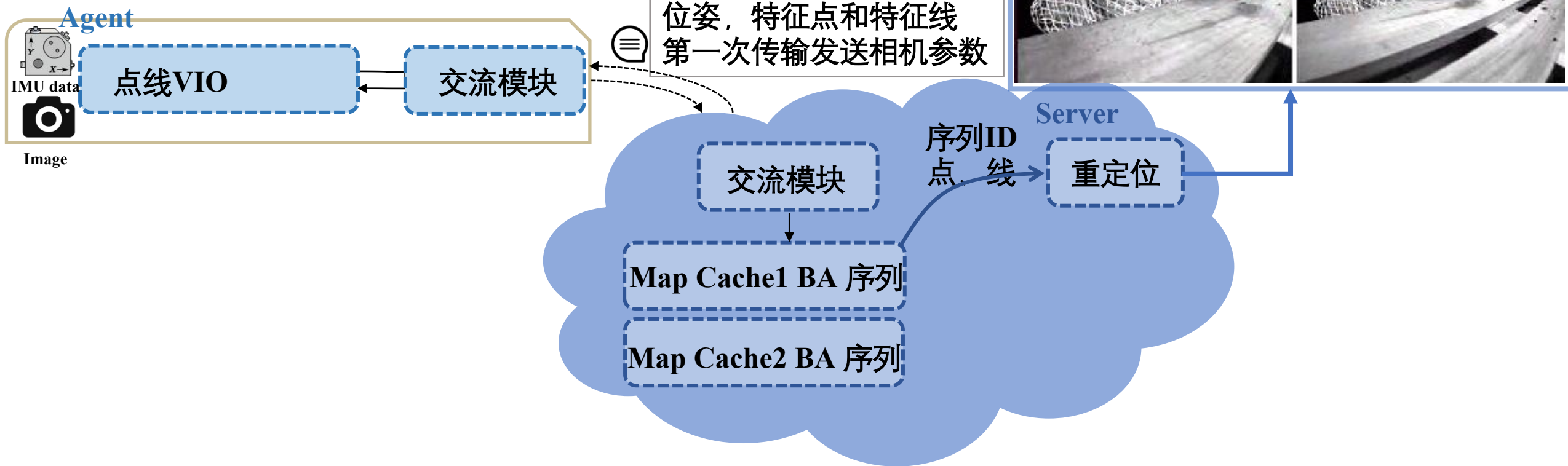


CoSLAM 包括一个两阶段异步优化方法，用于缓存子图（CSG）的优化和全局地图的优化。该方法确保对代理的实时响应，同时考虑了准确性和通信效率，显著提高了协同SLAM 的可扩展性

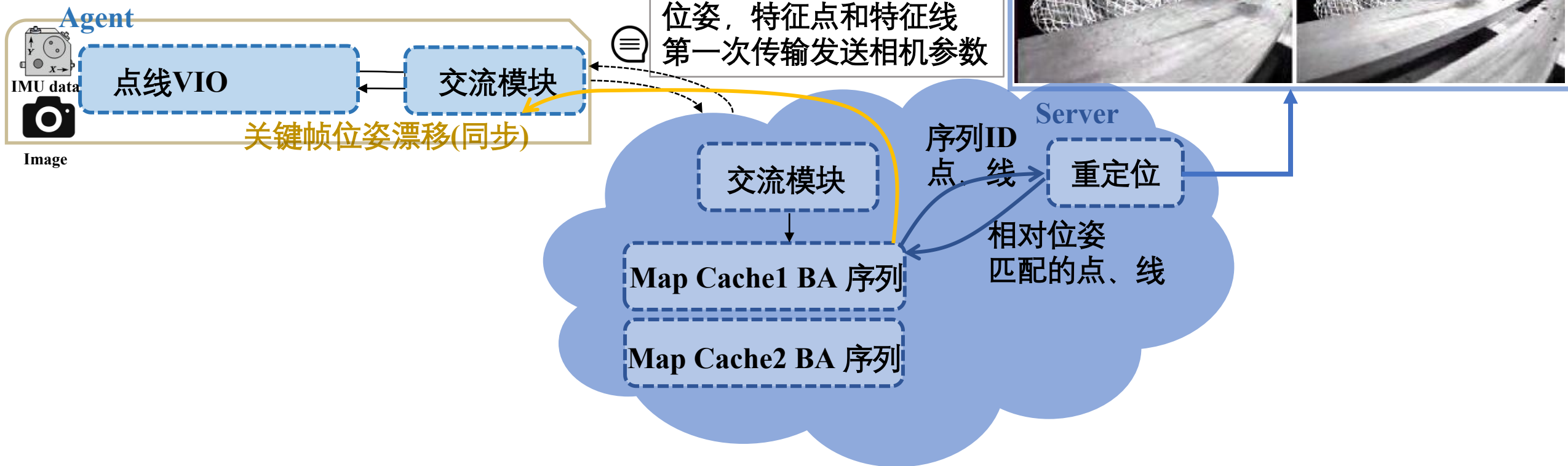


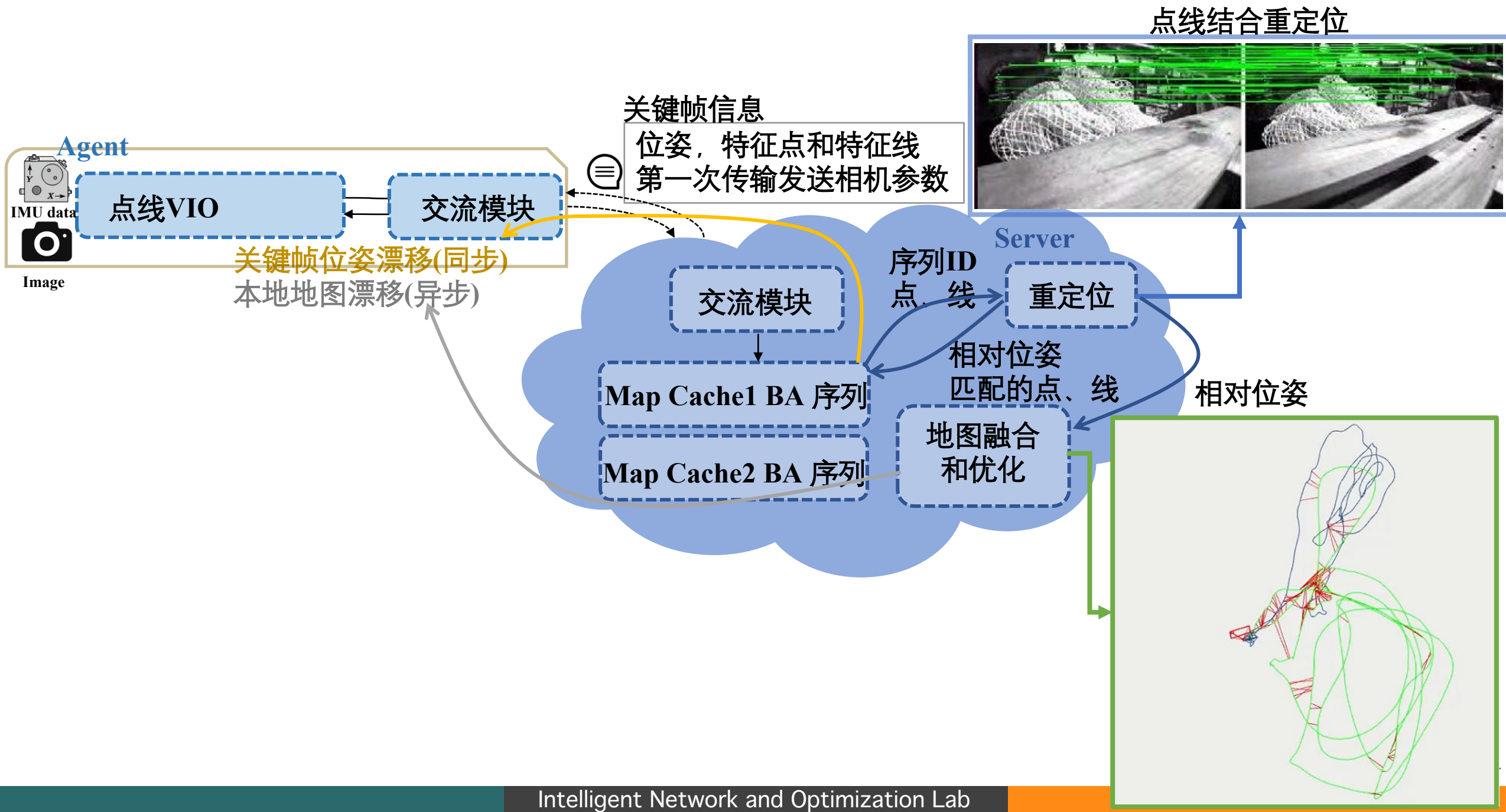


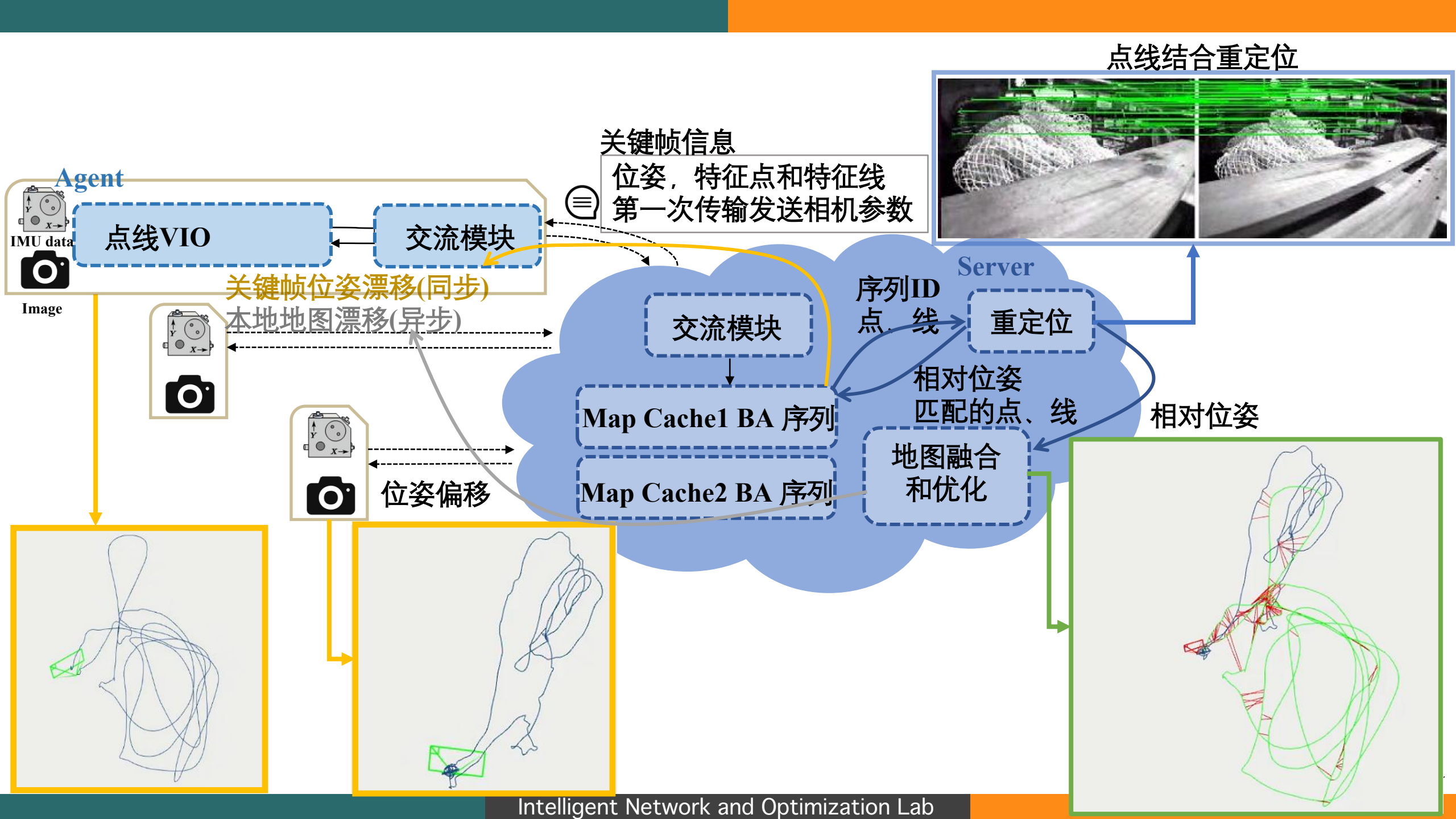
点线结合重定位



点线结合重定位







ColSLAM: A Versatile Collaborative SLAM System for Mobile Phones Using Point-Line Features and Map Caching

Demo Video

总结

- DroneMOT: 无人机多目标追踪系统
- GSLAMOT: 同步定位建图与多目标追踪系统
- SAHNet/RoCo: 多车协同感知系统
- MobiSketch: 基于点线面融合的语义3D建图
- CoISLAM: 多手机协同SLAM系统

谢谢, Q&A

ycw@ruc.edu.cn

主页: <http://www.yongcaiwan.com>
<http://yongcaiwan.github.io>