



全国汽车标准化技术委员会

National Technical Committee of Auto Standardization

# 《智能网联汽车感知训练数据集标准化需求研究报告》

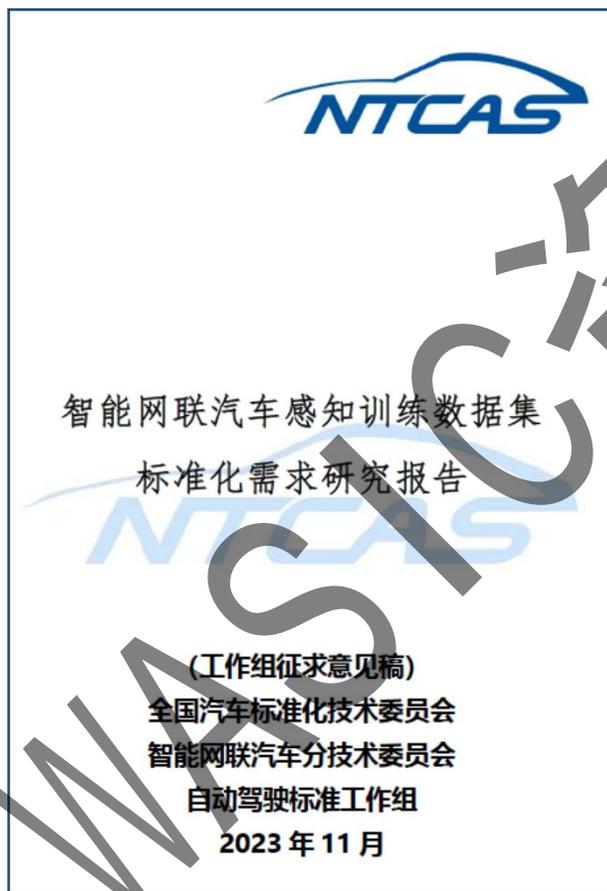
## 内容解读

中国汽车标准化研究院 智能网联部 华一丁

2023年11月27日

# 目录

## Contents



### 01 研究概况介绍

### 02 研究内容解读

### 03 标准输出建议

# 01 研究概况介绍

WASIC

## 1.1 研究背景

“ 现在，国际上**汽车行业**竞争很激烈，信息化、智能化等趋势不断发展，对我们来讲有危有机，危中有机，一定要把关键核心技术掌握在自己手里。

—— 2020年，**习近平总书记**在吉林省考察时来到中国一汽集团发表的讲话

“ 突破关键基础技术。开展复杂系统体系架构、**复杂环境感知**、……等基础前瞻技术研发。重点突破新型电子电气架构、**多源传感信息融合感知**、……等共性交叉技术。

—— 2020年，**中国国家发改委等11部门联合印发《智能汽车创新发展战略》**

- 随着技术水平的不断提升和相关产业的全面融合，我国智能网联汽车产业进入新的发展阶段。智能网联汽车标准体系建设第一阶段目标任务已圆满完成，初步构建起支撑驾驶辅助及低级别自动驾驶的标准体系。高精度环境感知是智能网联汽车技术创新体系中的关键基础技术之一，**具有中国特色感知数据训练集**又是高精度环境感知技术的重要组成部分，其质量一定程度上决定着感知能力的水平。

## 1.1 研究背景

- 为贯彻落实《国家标准化发展纲要》《国家车联网产业标准体系建设指南（智能网联汽车）》等文件要求，推进中国智能网联汽车标准体系建设，汽标委智能网联汽车分标委（SAC/TC114/SC34）秘书处启动首批**智能网联汽车（ICV）标准化领航项目研究**，其中包括**感知数据训练集**的标准化需求研究。



### 中国汽车标准化研究院是中汽中心直属技术机构

- > 中国唯一专业从事汽车标准化研究的国家级科研院所
- > 国家标准委、工信部授权的汽车标准化归口管理机构
- > 协助政府开展汽车国际标准法规协调的对口支撑单位

ICV 标准化领航项目 1-智能网联汽车 量子通信技术及其安全应用

ICV 标准化领航项目 2-基于先进通信技术的车辆网联功能与应用

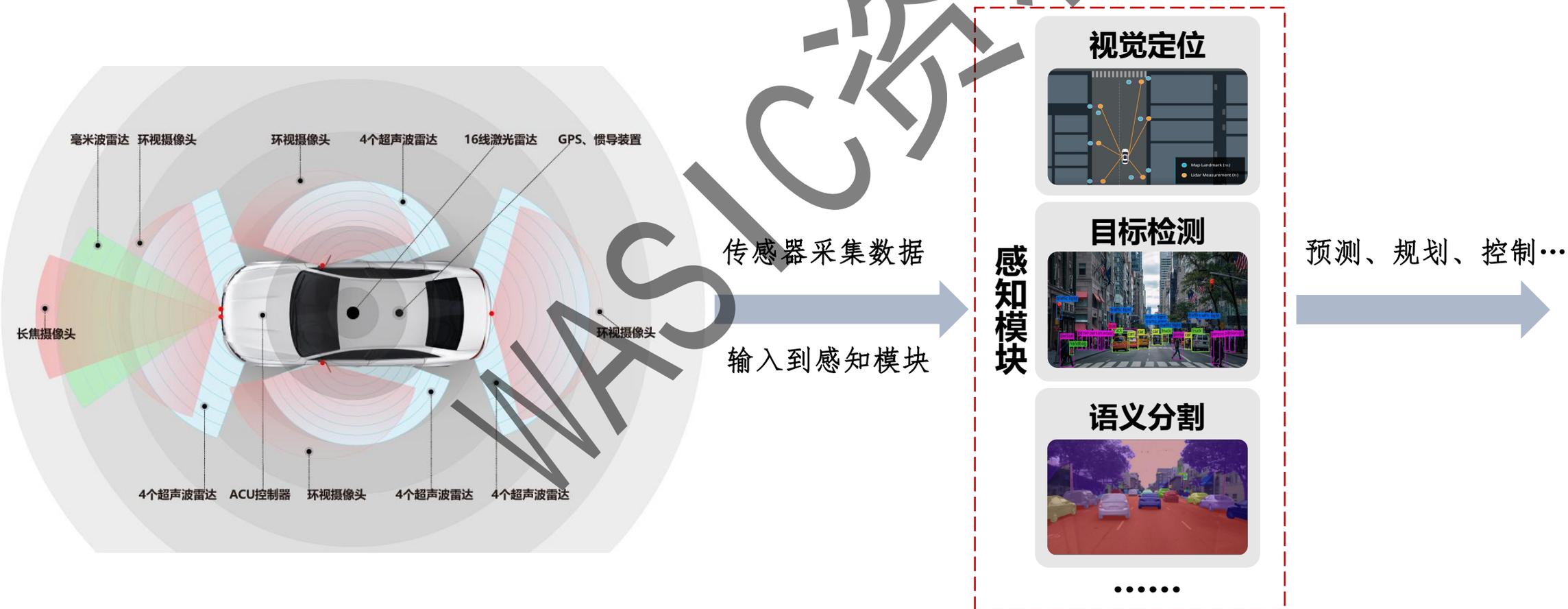
**ICV 标准化领航项目 3-智能网联汽车 感知数据训练集**

ICV 标准化领航项目 4-智能网联汽车 坐标系

ICV 标准化领航项目 5-智能网联汽车 云控平台

## 1.2 研究目的

- **感知训练数据集**是自动驾驶环境感知算法训练与实现的重要基础。为了实现高安全高可靠的自动驾驶，感知训练数据集要应对**多重挑战**：交通场景复杂性、天气多变性、光照挑战、传感器误差以及标注一致性等。
- **标准化**不仅是解决问题、提升感知算法鲁棒性的重要一环，更是引领自动驾驶技术蓬勃发展的关键步骤。



## 1.3 研究意义

- 构建具有中国特色感知数据训练集，并兼顾现有国外成熟感知数据训练集，将有效促进我国智能网联汽车感知算法能力提升，并为中国自主品牌“走出去”提供感知算法基础支撑。

### 为行业相关技术研发 提供基础支撑

1

ADAS/ADS技术方案不断成熟，产品驾驶自动化等级逐步提高，企业对更高精度的感知能力的需求增加，目前广泛采用的训练集往往是国外研究机构提供的，与中国的道路交通参与者的特征相差加大，存在较大的不适应性。

### 为行业管理提供 间接保障

2

政府对ICV的管理要求逐渐明晰，技术要求和测试方法相关标准加快发布实施，目前对于训练集测评手段及测试方法还处于前期研究阶段，标准化的测试方法是支撑标准实施的重要保障。

### 标准、高效、全面的 感知数据训练平台

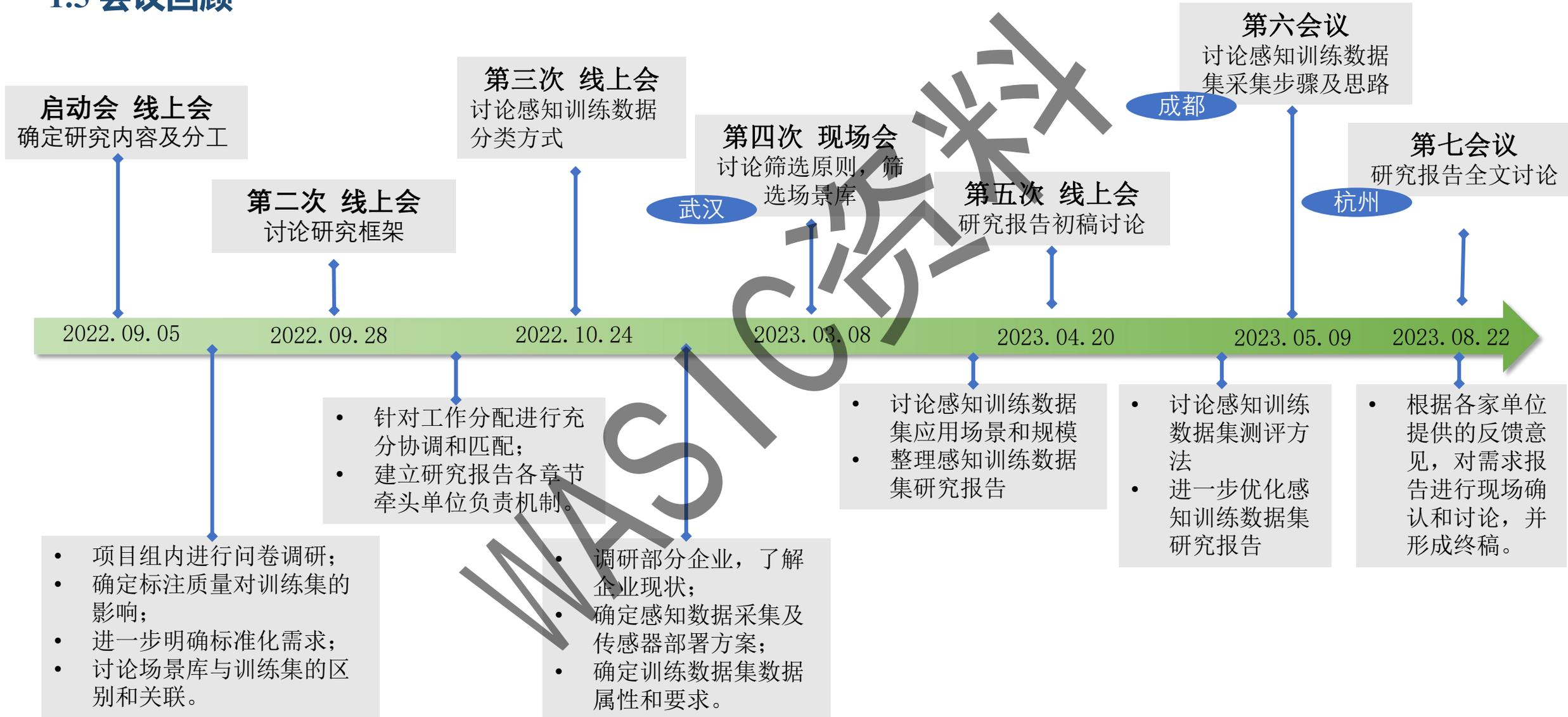
3

目前国内及国外存在诸多感知训练集，数据体量几千到几万帧不等，同时质量也难以保证，误标率、错标率难以支撑车企量产车算法，同时存在训练集覆盖度较为片面，缺乏统一及全面的统筹。

## 1.4 研究项目组成员

序号	单位名称	职责	序号	单位名称	职责
1	中国汽车技术研究中心有限公司	牵头	16	上海机动车检测认证技术研究中心有限公司	参与 单位
2	上海交通大学	参与 单位	17	小米汽车科技有限公司	
3	北京觉非科技有限公司		18	招商局检测车辆技术研究院有限公司	
4	北京赛目科技有限公司		19	福特汽车（中国）有限公司	
5	北京百度智行科技有限公司		20	上汽通用五菱汽车股份有限公司	
6	宁波吉利汽车研究开发有限公司		21	高通无线通信技术（中国）技术有限公司	
7	一汽解放汽车有限公司		22	西华大学	
8	上海临港绝影智能科技有限公司		23	武汉路特斯科技有限公司	
9	重庆长安汽车股份有限公司		24	江苏大学	
10	中国软件评测中心		25	天津大学	
11	北京智能车联产业创新中心有限公司		26	华为技术有限公司	
12	北京云测信息技术有限公司		27	东风悦享科技有限公司	
13	东风汽车集团有限公司技术中心		28	江铃汽车股份有限公司	
14	长城汽车股份有限公司		29	采埃孚商用系统有限公司	
15	泛亚汽车技术中心有限公司		30	高新兴科技集团股份有限公司	

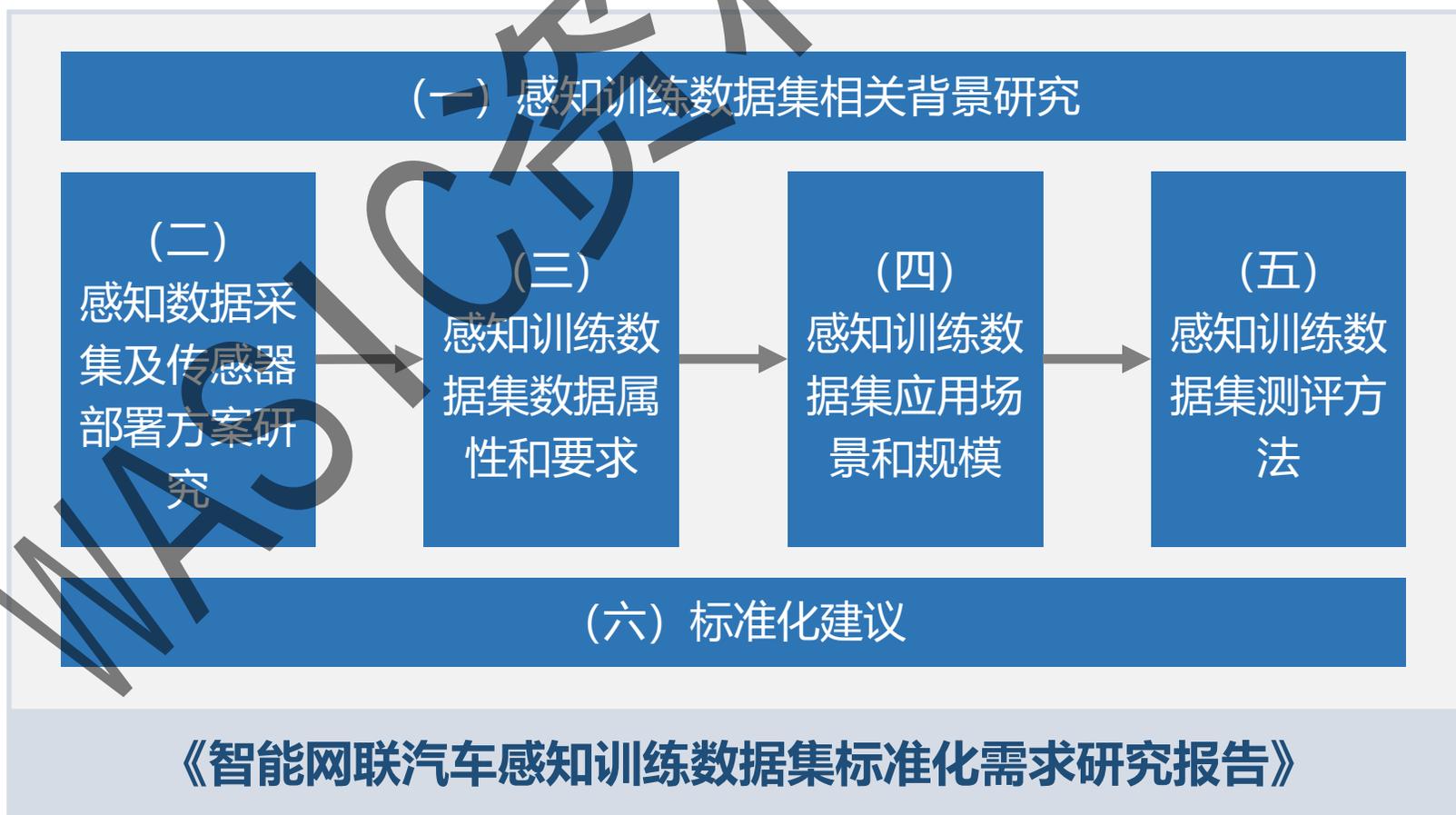
## 1.5 会议回顾



# 02 研究内容解读

WASIC

- 《智能网联汽车感知训练数据集标准化需求研究报告》（以下简称“本报告”）共6万余字，从感知训练数据集研究背景、感知数据采集及传感器部署方案、感知训练数据集应用场景和规模、感知训练数据集数据属性和要求、感知训练数据集测评方法以及标准化建议**六个方面**进行梳理和总结。



### (一) 感知训练数据集研究背景

- 根据采集途径不同，感知训练数据集可分为**车端**和**路端**两类。本报告调研了现有的车路两端感知训练数据集和常见的数据标注方式。

#### I. 车端感知训练数据集现状

- 车端数据集复杂多样，不仅包含驾驶数据集，还有交通标志数据集、行人数据集、车道线检测数据集等。
- 车端数据使用数据采集车收集，传感器部署方式多变。



高德采集车示意图



百度ApolloScape数据采集车示意图

## (一) 感知训练数据集研究背景

### I. 车端感知训练数据集现状

数据集	发布时间	地点	使用设备	场景	标注类别	数据规模
nuScenes	2019	波士顿和新加坡	1个32线雷达 6个摄像机	市区、住宅、 郊区和工业区	图像23类， 点云32类	1.4M高清图像，4万帧点云
KITTI	2012	德国卡尔斯鲁厄	1个64线雷达 4台摄像机	市区、乡村和高速	9类	29GB，15000个数据文件，超过200k 3D标注物体图像
Argoverse	2019	皮茨堡和迈阿密	2个32线雷达 9个摄像机	市区	15类	跟踪113个场景，预测32万条轨迹，每个轨迹5秒
ONCE	2021	中国多个城市	1个激光雷达 7个摄像机	市中心、郊区、 隧道、高速公路、 桥梁等	9类	1.6万个场景，41万个3D框和76万个2D边框
Waymo	2019	美国六个城市	5个64线雷达 5个摄像头	市区、郊区等	23类	检测2030个场景，预测10万个场景，113k个3D轨迹和160k个2D轨迹
ApolloScape	2018	中国十个城市	2个250线雷达 6个摄像机	市区，乡村和高速	35类	144K+张图像，70K帧3D实例标记，1000km行驶轨迹

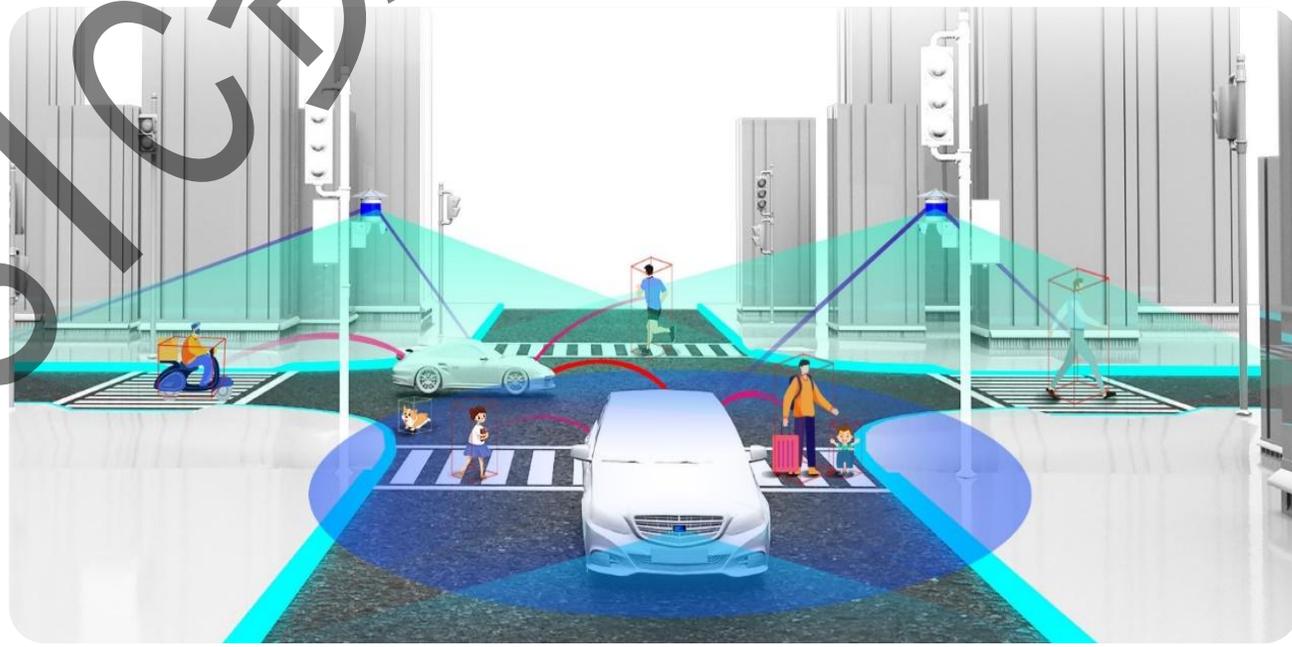
### (一) 感知训练数据集研究背景

#### II. 路端感知训练数据集现状

- 路端数据采集最主要的部署场景为城市中的交通路口，城市交通路口人员密集，障碍物众多，对于路端辅助感知有迫切需求，现有的路端感知数据集仍较少。
- 路端感知设备大多安装在高3-5m的设备杆上，或现有交通设施上。



路端感知设备示意图



车路协同感知示意图

### (一) 感知训练数据集研究背景

### II. 路端感知训练数据集现状

数据集	发布时间	数据量	数据类别	其他
DAIR-V2X	2022	71254帧图像数据 71254帧点云数据	10类目标物	包含障碍物遮挡、截断等信息
Rope3D	2022	50009帧图像数据	9类目标物	
IPS300+	2021	14198帧图像数据 14198帧点云数据	7类目标物	
BAAI-VANJEE Roadside Dataset	2021	5000帧图像数据 2500帧点云数据	12类目标物	
V2XSim	2021	三个Carla小镇 500G以上数据	23类目标物	大型虚拟数据集
A9-Dataset	2022	1098帧数据	9类目标物	包含大雪、大雾、交通事故特殊场景

### (一) 感知训练数据集研究背景

### III. 数据集标注能力现状

- 数据集标注是指将数据集中的原始数据（图片、点云、文本等）进行人工或自动处理，将其转换成计算机能够理解、能够被算法处理的标准形式。数据标注的质量与能力直接影响到机器学习算法的性能。
- 本报告调研了**通用**的数据集标注方式，以及针对**不同任务**的标注流程。

标注方式	通用标注方式			根据自动驾驶任务划分		
	人工标注	半自动标注	自动标注	传感器融合标注	路径与轨迹标注	场景理解标注
方式介绍	人工清洗标注	自动初步标签，人工审阅	纯算法标注，无需人工	多个传感器融合	目标跟踪、轨迹预测任务	标注场景相关信息
标注步骤	<ul style="list-style-type: none"> <li>a. 标注需求理解</li> <li>b. 标注准备</li> <li>c. 标注过程</li> <li>d. 标注结果</li> </ul>	<ul style="list-style-type: none"> <li>a. 算法预处理和预标注</li> <li>b. 人工修正和补充</li> <li>c. 标注结果</li> </ul>	<ul style="list-style-type: none"> <li>a. 算法处理和标注</li> <li>b. 标注结果</li> </ul>	<ul style="list-style-type: none"> <li>a. 数据同步</li> <li>b. 数据对齐</li> <li>c. 标注物体</li> <li>d. 标注关联</li> <li>e. 质量检查</li> </ul>	<ul style="list-style-type: none"> <li>a. 数据收集</li> <li>b. 检测物体标注</li> <li>c. 跟踪物体标注</li> <li>d. 轨迹生成</li> <li>e. 轨迹标注</li> <li>f. 质量检查</li> </ul>	<ul style="list-style-type: none"> <li>a. 道路属性标注</li> <li>b. 环境属性表述</li> <li>c. 上下文信息标注</li> <li>d. 场景分类标注</li> </ul>

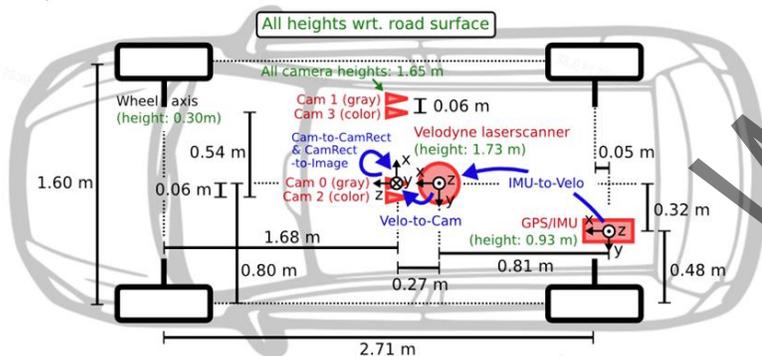
## (二) 感知数据采集及传感器部署方案研究

- 本报告分别对**车端**和**路端**两类感知数据采集及传感器部署方案进行了研究。

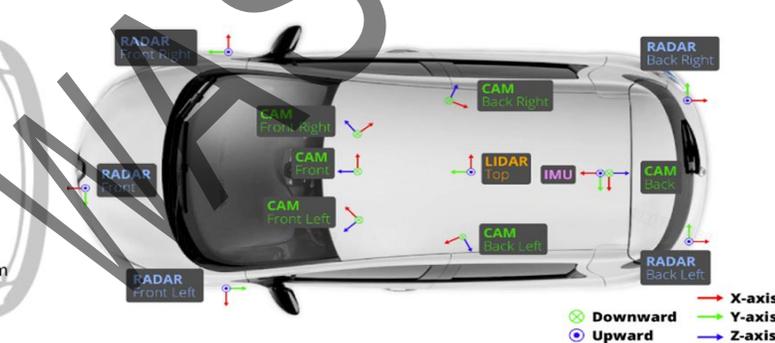
### I. 车端感知数据采集

采集车的种类及型号	传感器的选型	传感器配置数量及部署安装方案	采集系统配置	其他硬件系统
为最大化满足业务需求，选车需要既要考虑通用需求，也要关注特殊需要。一般包括车身外形与尺寸、车内空间、车顶行李架、能源类型及油标号等等十多项参数。	智能驾驶汽车环境感知传感器主要有激光雷达、单/双/三目摄像机、鱼眼摄像机、环视摄像机以及毫米波雷达等，选型时需要考虑通用水平和性能水平。	一般采用5-8个摄像机并根据功能，分别装在车顶、后视镜下方、前侧翼子板、前后车标、后备箱等。常采用1-3激光雷达个布置在采集车正上方离地2m-3m处。	主要包括外部感知设备、与外部设备相连接的交换机、同步盒、电源模块，数据融合单元，数据采集与存储系统，最终上传到工控机或云端服务器进行存储。	主要为组合导航。

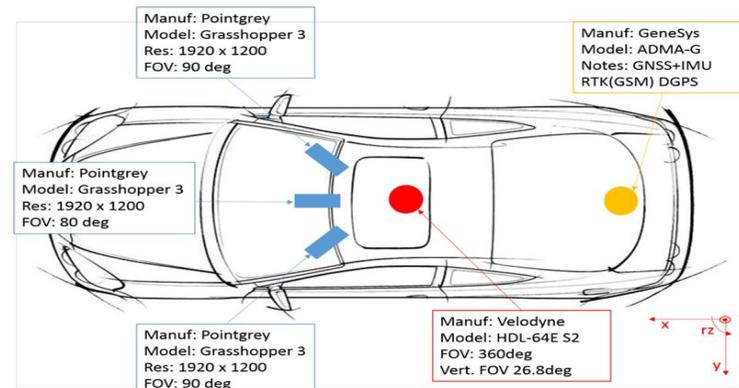
- 目前较为成熟的车端感知数据采集车传感器布置方案



KITTI数据集采集车传感器布置



nuScenes数据集采集车传感器布置



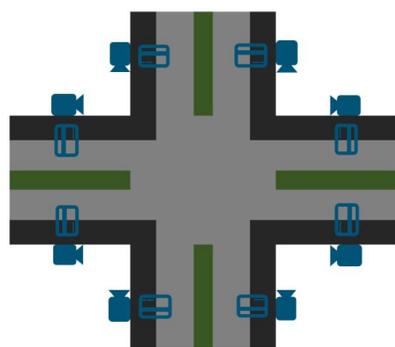
H3D数据集采集车传感器布置

## (二) 感知数据采集及传感器部署方案研究

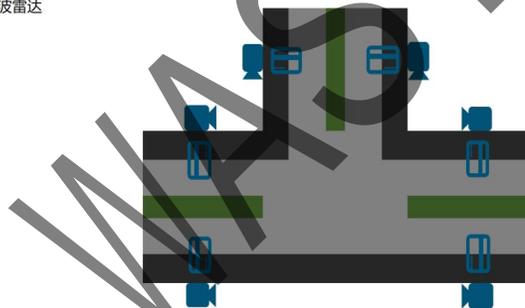
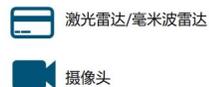
### II. 路端感知数据采集

传感器选型	传感器配置数量及部署安装方案	采集软件配置	其他硬件系统
路段环境感知传感器主要有激光雷达、毫米波雷达、RGB相机、红外相机、事件触发相机等选型时需要考虑通用水平和性能水平。	路端传感器主要部署场景包括城市路口、城市道路、高速公路等，对于不同的道路场景，需进行针对性的规划设计，可以划分为十字路口，T型路口，城市道路。	主要包括数据获取模块，目标融合模块，数据处理模块，孪生可视化平台模块。	主要为边缘计算单元。

#### ● 目前较为成熟的路端感知数据采集传感器布置方案



十字路口有绿化带布设方案



T型路口有绿化带布设方案



城市道路布设方案

## (三) 感知训练数据集数据属性和要求

- 本报告调研了针对不同感知任务的数据集标注范围、标注类别和标注属性。

### I. 视觉图像算法数据属性和要求

图像2D目标检测算法数据集标注范围（节选）

类别	描述
Vehicle (车辆)	此类别包含场景中出现的所有车辆类别，包括汽车，自行车，摩托车等子类。
Human (人)	此类别包含场景中检测到的所有人，包括行人，骑手，移动单元，人脸等子类别。
Traffic sign (交通标志)	此类别包含所有面向采集车方向的交通标志（包含交通灯），包括我们不需要标注标志的背面。
.....	.....

图像2D标注目标类别属性定义（节选）

属性	值	描述
Occluded (遮挡)	0: 0% (完全可见) 1: 1-50% (部分遮挡) 2: 51-100% (大部分遮挡) 3: 未知	此属性默认值为0，只有出现部分遮挡或全部遮挡时，将其更改为相应数值。
Truncated (截断)	0: 0% (完全可见) 1: 1-50% (部分可见) 2: 51-100% (大部分被截断)	截断指的是物体离开图片边缘的情况。适用于当车辆，人，交通标志，交通灯的部分或全部被截断的情况。默认值为0。
.....	.....	.....

### II. 激光点云算法数据属性和要求

点云3D目标检测标注范围（节选）

类别名称	类别说明
car	包含轿车、SUV、MPV、皮卡车、五菱宏光类小面，四轮快递车等。
bus	包含公交汽车、校车、小巴、大金杯、C型房车等。
truck	大型集卡、大型货车、大型半挂车、大型厢式货车、中型厢式货车、油罐车、垃圾车、机动农用三轮车等。
.....	.....

3D目标检测目标属性说明

属性	含义	值	值描述
truncated	是否截断：目标物是否在边缘且未完全在画面中	0	目标物完整
		1	目标物不完整（出现在固态雷达点云边缘）
occluded	是否被遮挡：目标物是否被其他物体遮挡	0	完全未被遮挡
		1	小部分被遮挡（不超过一半）
		2	大部分被遮挡
abnormal	数据是否异常	3	未知（不知道是否被遮挡）
		0	默认0，无异常
		1	残影（拼接不完善）

### (三) 感知训练数据集数据属性和要求

#### III. 训练数据标注精确度行业普遍情况

精确度指标	含义
抽检比例	待验收的已标注样本占总的已标注样本的比例
准确率	合格的已标注样本数占总的已标注数的比例
贴合度	作业人员标注的点、线、框与原始数据中的点、线、框的真值之间的重合程度
航向角	在传感器局部坐标系下，标注框的方向与局部坐标系的航向参考轴之间的夹角
连续帧的最短帧间隔时长	最短帧间隔时长指的是传感器连续帧之间的时间间隔
速度的标注以及最大误差	标注过程中允许的物体速度值与其真值之间差值的最大绝对值
图像标注的像素误差	2d标注框与图像中物体的真值框之间的像素差值

#### IV. 训练集数据标注格式

图像标注数据导出格式

标注项	导出格式
2D框	宜采用txt/json
语义分割	宜采用json/Mask
目标检测	宜采用xml/txt/json
目标追踪	宜采用xml/txt
车道线检测	宜采用json

点云标注数据导出格式

标注项	导出格式
3D框	宜采用json
语义分割	宜采用pcd/coco/voc
目标检测	宜采用pcd/bin/npz
目标追踪	宜采用xml/json
车道线检测	宜采用json

## (四) 感知训练数据集应用场景和规模

- 本报告调研了基于图像、激光点云以及数据融合算法对不同数据集的使用情况以及不同感知任务所需数据量级和人物模型所对应的场景维度。

### I. 国内外主流的感知算法模型以及训练数据集子集所需量级建议

国内外主流的感知算法模型以及训练数据集子集所需量级建议

算法类型	算法名称	数据集	所需量级建议
目标检测	Fast R-CNN、SSD	PASCAL VOC、MS-COCO	建议为达到95%以上的识别准确率和召回率，需要不少于100万张用于训练和验证的图像，每张图像平均具备8个以上的细分类别目标
目标跟踪	SORT、DeepSORT	MOT challenge	针对每一类目标，应不少于20万张用于训练和验证的图像（要求单目标时间连续5s以上），保证每帧中均含有有效目标。
语义分割	HRNet	Cityscapes、Semantic3D	针对图像语义分割算法需求，应不少于50万张用于训练和验证的图像，每张图像平均具备2个以上的细分类别目标，主要针对道路相关内容进行采集。
车道线检测	SCNN、RESA	Cityscapes、TuSimple	针对车道线检测，建议为达到99%以上的识别率与准确率，需要不少于50万张用于训练和验证的图像，尽可能覆盖多样性道路场景，如城市道路（包括复杂路口）。尽量覆盖各种车道线类型、道路情况、天气情况等。
车位线检测	DeepPS、DMPR-PS	ps2.0、PIL_PARK	为了实现精确的停车位检测，建议数据集数据量达到10万张以上，如果使用单个摄像机所得图片，建议数据量达到15万张以上。
图像点云融合检测	MV3D、PointFusion、Frustum-PointNet	KITTI、SUN-RGBD	在图像点云融合检测算法数据集应有图像点云的总数100w数量以上，目标数量大概的分布在5-50之间，保证平均每帧的目标大概在10个左右，覆盖城区高速快速路以及多个城市。
图像点云融合跟踪	DSM、mmMOT、EagerMOT	KITTI、NuScenes	建议数据集中总帧数应超过100w帧，目标跟踪的种类至少应包含车、行人类别。

## (四) 感知训练数据集应用场景和规模

### II. 模型任务所应用的场景维度

道路场景类型				交通参与者与标志		
公共道路	封闭道路	路面情况			天气与光照	
		交通参与者类型	交通参与者行为	交通标志物		
关键道路特征						
分叉路	铁架桥	隧道	匝道	环行路		

### (五) 感知训练数据集测评方法

- 由于目前各种公开感知训练数据集的侧重点不同，不同数据集经过训练后在同一感知算法进行测评时结果的差异性太大，导致无法实现在统一标准下进行相关测试及评价，因此，本报告中对评价指标和评价方法进行了调研和总结。

#### I. 感知数据质量评价指标

- 图像数据质量评价指标

- 1) 像素均值：图像像素的平均值；
$$u = \frac{1}{MN} \sum_{i=1}^M \sum_{j=1}^N F(i, j)$$
- 2) 图像标准差：图像像素灰度值相对于均值的离散程度；
$$\text{std} = \sqrt{\frac{1}{MN} \sum_{i=1}^M \sum_{j=1}^N (F(i, j) - u)^2}$$
- 3) 图像平均梯度：图像的清晰度；
$$\nabla G = \frac{1}{MN} \sum_{i=1}^M \sum_{j=1}^N \sqrt{\Delta x F(i, j)^2 + \Delta y F(i, j)^2}$$

- 点云数据质量评价指标

- 1) 点云密度：描述单位面积上激光雷达点的平均数量。
$$\rho = \frac{n - \sum_{i=1}^m n_i}{A - \sum_{i=1}^m A_i}$$
- 2) 高程精度：评价点云数据的高程与其真实的地面高程之间误差分布离散程度。
$$Z_\sigma = \pm \sqrt{\frac{\sum_{i=1}^n (Z_i - \bar{Z})^2}{n-1}}$$
- 3) 平面精度：评价点云数据的平面位置与其真实的地面位置之间误差分布的离散程度。
$$d_{XYmax} = \max \left[ \sqrt{(X_i - \hat{X}_i)^2 + (Y_i - \hat{Y}_i)^2} \right]$$

### (五) 感知训练数据集测评方法

#### II. 感知训练数据集评价方法

- 感知数据集数据类型及分布评价

评测数据集应做到场景分布多样性，包含标注类别多样，数据标注精度准确等维度进行评测。以下为不同分布维度下，道路场景数据集细类分布：

- 按天气、气候、时间分布、视角等分布
- 按交通场景分布
- 按交通参与者统计

- 数据集划分方式

- 留出法**：按照固定比例将数据集静态的划分为训练集、验证集、测试集；
- 留一法**：每次的测试集都只有一个样本，要进行m次训练和预测。这个方法用于训练的数据只比整体数据集少了一个样本，因此最接近原始样本的分布
- k折交叉验证**：一种动态验证的方式，这种方式可以降低数据划分带来的影响

- 评估指标

任务类型	指标
分类问题	准确率/差准率/召回率/F1分数 /ROC曲线/AUC
回归问题	平均绝对误差/均方误差/均方根 误差/均值平方对数误差/平均绝对百分比误差

- 根据调研结果，本报告提出了感知训练数据集标准化建议

### 标准内容

- 建议开展对车辆感知数据的数据标注内容、标注质量等具体方向的行业标准研制工作

### 标准支撑

- 目前开展了“自动驾驶系统测试场景数据采集和分析标准需求研究”项目，建议参考该项目的研究进程，进一步明确相关的标准化需求情况。

### 注意事项

- 数据标注现有即将正式实施的推荐性国家标准《GB/T 42755-2023 人工智能 面向机器学习的数据标注规程》，建议进行区分。

# 03 标准输出建议

WASIC

### 1 充分考虑现有感知训练数据集的差异性，有针对性地开展通用技术标准化

- 标准化感知训练数据集为行业提供标准、高效、全面的感知数据训练平台，为相关技术研发提供基础支撑；
- 目前，因感知数据集领域的相关标准较少，尤其是开发车辆感知功能所用的数据集根据目标功能、技术路线等的不同，在行业内普遍存在一定的差异，标准化过程应充分考虑传感器类型、训练目的及应用场景等内容。

### 2 进一步关注自动驾驶大模型等新型人工智能技术对训练数据集标准化的需求

- 实现可信AI，数据的设计、改进和质量评估是关键。以数据为中心的人工智能将重点转移到训练数据的治理和增强，高质量的训练数据集、完备的数据应用策略将会更好的服务于模型的开发与应用；
- 价值对齐是AI安全的核心议题。大模型的能力和行为跟人类的价值、真实意图和伦理原则需保持一致，为确保人类与人工智能协作过程中的安全与信任，需重点关注相关自动驾驶感知训练数据集的标准化研究。

### 3 重视行业相关标准确实问题，加速推动自动驾驶感知技术落地应用

- 国内针对车辆感知数据集数据标注的标准仍处于空白状态，建议重点针对车辆感知数据的数据标注内容、标注质量等具体指标的开展行业标准研制工作。



全国汽车标准化技术委员会  
National Technical Committee of Auto Standardization

请各位领导、专家  
批评、指正！

汽标委智能网联汽车分委会 华一丁

- 022-84379865
- huayiding@catarc.ac.cn
- 18622766087