**附件1**

# 浙江工程师学院（浙江大学工程师学院）
# 同行专家业内评价意见书

姓名：　　　　　　　　　杨一帆

学号：　　　　　　　22260198

申报工程师职称专业类别（领域）：　　　交通运输

浙江工程师学院（浙江大学工程师学院）制

2025年05月21日

# 填表说明

一、本报告中相关的技术或数据如涉及知识产权保护、军工项目保密等内容，请作脱密处理。

二、请用宋体小四字号撰写本报告，可另行附页或增加页数，A4纸双面打印。

三、表中所涉及的签名都必须用蓝、黑色墨水笔，亲笔签名或签字章，不可以打印代替。

四、同行专家业内评价意见书编号由工程师学院填写，编号规则为：年份4位＋申报工程师职称专业类别(领域)4位+流水号3位，共11位。

# 一、个人申报

**（一）基本情况【围绕《浙江工程师学院（浙江大学工程师学院）工程类专业学位研究生工程师职称评审参考指标》，结合该专业类别（领域）工程师职称评审相关标准，举例说明】**

## 1. 对本专业基础理论知识和专业技术知识掌握情况（不少于200字）

本人系统掌握了智能交通系统与自动驾驶技术领域的核心理论与关键技术，具备扎实的自动控制、感知融合、路径规划及决策制定等专业基础知识。在深入研究复杂交通环境下的多智能体系统行为与自主决策规划过程中，熟悉并掌握了强化学习、深度学习等人工智能算法在智能交通中的应用，能够结合不同场景需求进行模型设计与算法优化。通过对多源传感器数据的融合处理，构建了高精度的环境感知模型，为自动驾驶系统的稳定运行提供了有力支撑。

在专业技术方面，掌握了以ROS为核心的机器人开发框架，能够熟练使用Python、MATLAB等工具进行算法实现、系统仿真与调试优化。深入研究了面向多智能体系统的路径规划与协同控制技术，提出的强化学习与规则融合的决策算法在仿真与实验中表现出良好的泛化能力和环境适应性。同时具备较强的工程实践能力，能够结合SLAM、自主导航、目标检测等技术，实现自动驾驶车辆在复杂场景中的稳定运行。通过参与多项项目开发与企业实践，积累了丰富的工程经验，并具备将先进算法应用于物流自动化、无信号控制十字路口等实际场景的能力，展现出良好的技术创新与工程转化能力。

## 2. 工程实践的经历（不少于200字）

在"先进多智能体决策规划技术研究"项目中，我承担了项目负责人及核心研发成员的角色，深入参与了从理论研究到系统开发的全过程。在工程实践方面，我主导了多智能体协同算法在仿真与实际环境中的部署与测试，搭建了基于ROS和Gazebo的仿真平台，对路径规划与交互决策算法进行了多轮仿真评估与参数优化。在实验过程中，我针对阿克曼车辆的非线性运动学特性，开发并验证了适用于多智能体系统的优化路径规划算法，提升了机器人在复杂场景中的运动效率与稳定性。

此外，我参与了实际应用测试的准备工作，包括传感器调试、环境建模与数据采集，并协同团队构建了一个面向物流配送和搜索救援场景的实验测试环境。在测试过程中，我负责系统集成与调试，确保多智能体之间能够进行高效的信息交换与协同控制。通过与实际需求对接，我积累了宝贵的系统工程经验，并在应对数据噪声、通信延迟及环境动态变化等现实问题中不断调整优化策略，提升了项目的工程实用性和技术可靠性。这一系列实践经历不仅加深了我对多智能体系统工程实现的理解，也锻炼了我解决实际问题和推动技术落地的能力。

## 3. 在实际工作中综合运用所学知识解决复杂工程问题的案例（不少于1000字）

在"先进多智能体决策规划技术研究"项目的工程实践过程中，我带领团队针对多智能体系统在复杂环境中的协同决策与路径规划难题，进行了深入的理论探索与工程实现。该项目旨在解决多智能体系统在现实场景中执行任务时的稳定性、鲁棒性和协同效率问题，尤其在物流配送和搜索救援等典型应用中，系统需要面对环境动态性强、智能体交互复杂等工程挑战。整个项目过程贯穿了从理论建模、算法设计、系统开发到现场测试的完整闭环，我也在实践中综合运用了自动控制、车辆动力学、人工智能、优化理论、计算机视觉和嵌入式系统开发等多个专业领域的知识，推动了项目的有效实施。

我们从多智能体系统的建模与仿真开始着手。在理论研究阶段，我结合车辆动力学与非线性控制理论，针对阿克曼类转向结构的机器人建立了高保真的运动学模型。由于这种非线性结构在常规规划算法中难以精确控制，我开发了一种基于非线性模型预测控制的路径生成器，能够在路径生成阶段就考虑车辆的可控性与实际运行约束，确保轨迹的可行性与光滑性。在

路径生成过程中，为避免路径间冲突以及时间窗重叠的问题，我引入了多智能体间的时序协调机制，实现了任务分配与路径优化的联动规划。这一方法不仅提升了多车系统的协同效率，还有效缓解了交通瓶颈区域的拥堵情况。

在协同决策方面，我设计了一套融合博弈理论和强化学习的分布式策略。系统根据任务紧急程度、路径代价与智能体资源状态，为每一智能体分配最优的任务与路径。在这一过程中，我充分运用了在人工智能与最优化方法中所学的知识，将状态空间建模、Q-learning学习策略和冲突协调机制结合起来，构建了一个可以自适应动态环境变化的任务调度系统。在实际运行中，系统能够根据环境变化实时调整决策策略，实现多车间的高效协作。为确保系统的实时性和稳定性，我还搭建了基于ROS的通信架构，并优化了节点间的数据传输机制。通过发布-订阅模型，智能体间可实时共享状态、路径及感知数据，从而实现了分布式系统的高效协同。

感知系统是保障系统准确性的基础。在这一模块中，我采用了多传感器融合方法，将激光雷达、摄像头、GPS与IMU数据进行融合处理，构建了高精度的环境感知模型。通过扩展卡尔曼滤波器处理位姿信息，我实现了对移动障碍物的动态检测和跟踪。同时，结合深度学习技术，我设计了基于YOLO的目标识别模型，用于识别特定类型障碍物和任务目标，为路径规划和任务决策提供支持。在栅格地图的基础上，我进一步引入代价地图机制，使路径规划算法在评估可行路径时不仅考虑避障需求，还能参考实时交通密度、通行舒适性等参数，提升了路径生成的合理性。

整个系统的仿真验证在Gazebo平台上完成。我带领团队构建了多个仿真场景，包括复杂交叉路口、动态障碍物区域及多任务密集区域等。在这些环境中测试的过程中，我对算法的性能进行了全面评估，分析其在任务分配正确率、路径规划时间、交通冲突次数等指标上的表现。仿真结果显示，我们提出的协同算法在任务完成时间上比传统方法提高了约28%，在多智能体冲突处理能力上也有明显优势。这为后续的实际部署奠定了良好基础。

在实际测试阶段，我带领团队将系统部署到实验物流仓储环境中。我们使用Jetson嵌入式平台进行控制系统部署，并在多台具备激光雷达和摄像头的移动平台上安装自主导航系统。在现场测试中，我们进行了任务分配、路径规划、交互避障等多个场景的测试验证。在一次典型的任务测试中，三台移动机器人需要分别完成货物搬运、障碍物排查与物资配送任务。系统通过动态感知与任务调度，能够快速对环境变化做出响应，实现车辆之间的任务协商与路径重规划，成功完成任务执行。测试过程中系统运行稳定，未出现通信中断或路径冲突问题，验证了算法在真实环境中的适应性与实用性。

在整个项目过程中，我们也遇到了诸多问题。例如，系统在面对突发障碍物时路径重规划响应较慢，造成部分任务延误；传感器在强光干扰或低光环境下数据质量下降，影响目标识别准确性。对此，我主导优化了路径重规划算法，引入快速随机树（RRT）算法进行局部路径更新，并在传感器处理模块中增加滤波与补偿机制，提升系统鲁棒性。我们还与企业合作建立了数据采集机制，引入实际场景中的样本数据用于算法优化与模型训练，进一步提升了系统在工业环境中的适应能力。

通过这一工程项目，我不仅将理论知识有效转化为工程解决方案，还深刻理解了多智能体系统在实际部署中需要面对的挑战与技术细节。同时，这一经历也显著提升了我的系统集成能力、团队协调能力与工程创新思维。在理论与实践不断融合的过程中，我逐步构建起了从模型设计到系统实现再到实际验证的完整工程能力框架。这一项目的成功实施不仅推动了多智能体协同规划技术的工程落地，也为我未来在智能系统、自动驾驶与智能交通等领域的研究与实践奠定了坚实基础。

（二）取得的业绩（代表作）【限填3项，须提交证明原件（包括发表的论文、出版的著作、专利证书、获奖证书、科技项目立项文件或合同、企业证明等）供核实，并提供复印件一份】

1.
公开成果代表作【论文发表、专利成果、软件著作权、标准规范与行业工法制定、著作编写、科技成果获奖、学位论文等】

| 成果名称 | 成果类别<br>［含论文、授权专利（含发明专利申请）、软件著作权、标准、工法、著作、获奖、学位论文等］ | 发表时间/授权或申请时间等 | 刊物名称/专利授权或申请号等 | 本人排名/总人数 | 备注 |
|---|---|---|---|---|---|
| Intelligent Hybrid Decision-Making for High-Speed Autonomous Driving Scenarios | 会议论文 | 2024年12月10日 | The COTA International Conference of Transportation Professionals | 1/6 | |
| 一种多智能体高效协同路径规划方法 | 发明专利申请 | 2024年09月27日 | 申请号：202411359050.9 | 2/4 | 导师为一作，本人为二作，已进入实质审查阶段 |
| 一种基于分层搜索的多无人车构型保持协作运动规划方法 | 发明专利申请 | 2024年04月26日 | 申请号：202410512625.X | 3/4 | 导师为一作，已进入实质审查阶段 |

2. 其他代表作【主持或参与的课题研究项目、科技成果应用转化推广、企业技术难题解决方案、自主研发设计的产品或样机、技术报告、设计图纸、软课题研究报告、可行性研究报告、规划设计方案、施工或调试报告、工程实验、技术培训教材、推动行业发展中发挥的作用及取得的经济社会效益等】

| （三）在校期间课程、专业实践训练及学位论文相关情况 | |
|---|---|
| 课程成绩情况 | 按课程学分核算的平均成绩： 84 分 |
| 专业实践训练时间及考核情况（具有三年及以上工作经历的不作要求） | 累计时间： 1.2 年（要求1年及以上）<br>考核成绩： 87 分 |
| **本人承诺** | |

个人声明：本人上述所填资料均为真实有效，如有虚假，愿承担一切责任，特此声明！

申报人签名： 杨一帆

22260198

## 二、日常表现考核评价及申报材料审核公示结果

| 日常表现考核评价 | 非定向生由德育导师考核评价、定向生由所在工作单位考核评价：<br><br>☑优秀　☐良好　☐合格　☐不合格<br><br>德育导师/定向生所在工作单位分管领导签字（公章）：　　2025年5月21日 |
|---|---|
| 申报材料审核公示 | 根据评审条件，工程师学院已对申报人员进行材料审核（学位课程成绩、专业实践训练时间及考核、学位论文、代表作等情况），并将符合要求的申报材料在学院网站公示不少于5个工作日，具体公示结果如下：<br><br>☐通过　　☐不通过（具体原因：　　　　　　　　）<br>工程师学院教学管理办公室审核签字（公章）：　　　　年　月　日 |

# 浙 江 大 学 研 究 生 院

## 攻读硕士学位研究生成绩表

| 学号：22260198 | 姓名：杨一帆 | | 性别：男 | 学院：工程师学院 | | 专业：交通运输 | | | 学制：2.5年 |
|---|---|---|---|---|---|---|---|---|---|
| 毕业时最低应获：24.0学分 | | 已获得：26.0学分 | | | | 入学年月：2022-09 | | 毕业年月： | |
| 学位证书号： | | | 毕业证书号： | | | | 授予学位： | | |

| 学习时间 | 课程名称 | 备注 | 学分 | 成绩 | 课程性质 | 学习时间 | 课程名称 | 备注 | 学分 | 成绩 | 课程性质 |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 2022-2023学年秋季学期 | 工程技术创新前沿 | | 1.5 | 90 | 专业学位课 | 2022-2023学年秋冬学期 | 研究生英语 | | 2.0 | 86 | 公共学位课 |
| 2022-2023学年秋季学期 | 工程伦理 | | 2.0 | 85 | 公共学位课 | 2022-2023学年春季学期 | 科技创新案例探讨与实战 | | 2.0 | 85 | 专业选修课 |
| 2022-2023学年冬季学期 | 科学与工程计算方法 | | 2.0 | 81 | 跨专业课 | 2022-2023学年春季学期 | 自然辩证法概论 | | 1.0 | 82 | 公共学位课 |
| 2022-2023学年秋冬学期 | 研究生论文写作指导 | | 1.0 | 82 | 专业学位课 | 2022-2023学年春季学期 | 研究生英语基础技能 | | 1.0 | 68 | 公共学位课 |
| 2022-2023学年秋冬学期 | 高阶工程认知实践 | | 3.0 | 90 | 专业学位课 | 2022-2023学年春夏学期 | 优化算法 | | 3.0 | 78 | 专业选修课 |
| 2022-2023学年冬季学期 | 新时代中国特色社会主义理论与实践 | | 2.0 | 90 | 公共学位课 | 2022-2023学年夏季学期 | 智能交通系统与实践应用 | | 2.0 | 91 | 专业学位课 |
| 2022-2023学年冬季学期 | 产业技术发展前沿 | | 1.5 | 80 | 专业学位课 | | 硕士生读书报告 | | 2.0 | 通过 | |
| | | | | | | | | | | | |

说明：1. 研究生课程按三种方法计分：百分制，两级制（通过、不通过），五级制（优、良、中、

及格、不及格）。

2. 备注中"*"表示重修课程。

学院成绩校核章：

成绩校核人：张梦依

打印日期：2025-06-03

**Publication Certificate**

# The 25th COTA International Conference of Transportation Professionals

## (CICTP2025)

Issue Date:
February 24, 2025

This is to certify that the following paper:

**Title**

**Intelligent Hybrid Decision-Making for High-Speed Autonomous Driving Scenarios**

**Author(s)**

**Yifan Yang, Yuchen Wu, Gang Xu, Yong Liu, Zhitao Zhang, Jian Yang**

has been accepted for publication in the 25th COTA International Conference of Transportation Professionals (CICTP2025), Guangzhou, China | July 22-25, 2025.

Presented By:

*Dr. Guohui Zhang, CICTP2025 General Chair*

---

https://libdb.zju.edu.cn/s/lib/libtb/show/492

浙江大学 图书馆
ZHEJIANG UNIVERSITY

字母: A B C D E F G H I J K L M N O P Q R S T U V W X Y Z

学科: 人文 社科 管理 经济 数学 物理/力学 化学/化工/材料 地球科学/天文学 生物 农业/食品 医药
电子电气 计算机/网络 机械 土木/建筑 环境/能源 交通/航空/航天 综合

类型: 图书 期刊/会议论文 文学/位论文 专利/标准/法规 科技报告 事实数据 据纸 图片 文摘 目录 多种类型

数据库名称 [检索]

中文资源
外文资源
试用资源
免费资源

>>版权公告<< >>图书馆数字资源校外访问指南<<

**ASCE**

美国土木工程师协会数据库 [收藏]

网址: http://ascelibrary.org
校内访问方式: IP认证
校外访问方式: WebVPN、CARSI
购买类型: 正式购买资源
主要学科: 土木 建筑

访问内容年限: 期刊: 创刊年-; 会议录: 1996-
主要文献类型: 期刊 会议录 图书

[评价打分] (您的评价将作为数字资源的评估依据)

美国土木工程师协会 (American Society of Civil Engineers, 简称ASCE) 成立于1852年，至今已有150多年的悠久历史，是历史最悠久的国家专业土木工程师学会。现在，ASCE已是全球土木工程领域的领导者；所属的会员来自159个国家超过13万的专业人员。为了融脂在工程师之间分享更多的信息，ASCE已和其他国家的65个土木工程学会达成了合作协议。

2016年，ASCE数据库包含35种专业期刊 (已有27种被SCI收录，其中8种ASCE期刊的影响因子在各自学科系列时序排名前十)，超过425卷会议录，以及各种图书和标准。ASCE于2004年推出在线会议录 (ASCE Online Proceedings)，收录ASCE召开或与其他知名学会合办的国际会议的文献。会议录于实际应用，为土木工程从业者和研究者提供的创新技术和前沿技术发现的全面而深入的研究信息。ASCE会议录是土木工程领域的核心资源，并且无法从其他途径取得。

ASCE出版物涵盖学科：Mechanics (工程力学)、Management (工程项目管理)、Structural (结构)、Construction (施工)、Environmental (环境)、Urban Planning (城市规划)、Geotechnical (地质技术)、Water Resources (水资源)、Hydraulic (水力)、Coastal and Ocean (海岸和海洋工程)、Aerospace (航空宇宙)、Materials (建筑材料)、Architectural (建筑设计)、Professional Issues (建筑师职业)、Energy (能源)、Transportation (交通运输)、Infrastructure (基础设施)、Computing in Civil Engineering (土木工程领域的计算机应用)。

图书馆订购了ASCE的35种期刊、会议录、Civil Engineering Magazine Archive、GEOSTRATA杂志。ASCE会议信息可通过 http://www.asce.org/conferences_events查找。

更新:

---

07-1512010150-117656.xlsx - LibreOffice Calc

文件(F) 编辑(E) 视图(V) 插入(I) 格式(O) 样式(Y) 工作表(S) 数据(D) 工具(T) 窗口(W) 帮助(H)

A2 | Source title

EI COMPENDEX SOURCE LIST: UPDATED MARCH 17, 2017

| | A | B | C | D | E | F | G | H |
|---|---|---|---|---|---|---|---|---|
| | Source title | Source type | ISBN13 | Publisher | Country | Subject 1 | Subject 2 | Subject 3 |
| 8364 | AUA Guidelines for Backfilling and Contact Grouting of Tunnels and Shafts | book | 9780784406342 | American Society of Civil | United States | Civil and Structural Eng | | |
| 8398 | Automated People Movers and Automated Transit Systems 2016: Innovation in a Rapidly Urbanizing World - Proceedings of the 15th Internat | proceeding | 9780784479797 | American Society of Civil | United States | Civil and Structural Eng | Transportation | |
| 8399 | Automated People Movers and Transit Systems 2011: From People Movers to Fully Automated Urban Mass Transit - Proceedings of the 13th A | proceeding | 9780784411933 | American Society of Civil | United States | Transportation | | |
| 8400 | Automated People Movers and Transit Systems 2013: Half a Century of Automated Transit - Past, Present, and Future - Proceedings of the 14t | proceeding | 9780784477823 | American Society of Civil | United States | Biomedical Engineering | Transportation | |
| 8425 | Aviation: A World of Growth - Proceedings of the 29th International Air Transport Conference, IATC 2007 | proceeding | 9780784409381 | American Society of Civil | United States | Transportation | | |
| 8436 | Baltimore Civil Engineering History | proceeding | 9780784407992 | American Society of Civil | United States | Civil and Structural Eng | Building and Construc | |
| 8596 | Breakwaters '99 First Coastal Symposium on Monitoring of Breakwaters | proceeding | 9780784405970 | American Society of Civil | United States | Civil and Structural Eng | | |
| 8602 | Bridge Safety and Reliability | book | 9780784404423 | American Society of Civil | United States | Engineering (all) | | |
| 8604 | Bridging the East and West: Theories and Practices of Transportation in the Asia Pacific - Selected Papers from the Proceedings of the 11th A | proceeding | 9780784479810 | American Society of Civil | United States | Transportation | | |
| 8605 | Bridging the Gap: Meeting the World's Water and Environmental Resources Challenges - Proceedings of the World Water and Environmen | proceeding | 9780784405697 | American Society of Civil | United States | Water Science and Tech | | |
| 8621 | Building a Sustainable Future - Proceedings of the 2009 Construction Research Congress | proceeding | 9780784410202 | American Society of Civil | United States | Civil and Structural Eng | Building and Construc | |
| 8680 | Carbonate Beaches 2000 | proceeding | 9780784406403 | American Society of Civil | United States | Geochemistry and Petr | | |
| 8715 | Case Studies in Optimal Design and Maintenance Planning of Civil Infrastructure Systems | proceeding | 9780784404201 | American Society of Civil | United States | Engineering (all) | | |
| 8795 | Chimney and Stack Inspection Guidelines | book | 9780784406939 | American Society of Civil | United States | Safety, Risk, Reliability | | |
| 8818 | CICTP 2012: Multimodal Transportation Systems - Convenient, Safe, Cost-Effective, Efficient - Proceedings of the 12th COTA | proceeding | 9780784412442 | American Society of Civil | United States | Transportation | | |
| 8819 | CICTP 2015 - Efficient, Safe, and Green Multimodal Transportation - Proceedings of the 15th COTA International Conference of Transportati | proceeding | 9780784479292 | American Society of Civil | United States | Transportation | | |
| 8820 | CICTP 2016 - Green and Multimodal Transportation and Logistics - Proceedings of the 16th COTA International Conference of Transportation | proceeding | 9780784479896 | American Society of Civil | United States | Transportation | | |
| 8861 | Civil Engineering and Urban Planning 2012 - Proceedings of the 2012 International Conference on Civil Engineering and Urban Planning | proceeding | 9780784412435 | American Society of Civil | United States | Civil and Structural Eng | Geography, Planning a | |
| 8863 | Civil Engineering Education Issues 2001 Proceedings of the Thirth National Congress | proceeding | 9780784405901 | American Society of Civil | United States | Civil and Structural Eng | Education | |
| 8902 | Coastal Dynamics 2005 - Proceedings of the Fifth Coastal Dynamics International Conference | proceeding | 9780784408551 | American Society of Civil | United States | Geotechnical Engineer | Civil and Structural Eng | |
| 8903 | Coastal Engineering 2000 - Proceedings of the 27th International Conference on Coastal Engineering, ICCE 2000 | proceeding | 9780784405499 | American Society of Civil | United States | Ocean Engineering | | |
| 8904 | Coastal Engineering Practice - Proceedings of the 2011 Conference on Coastal Engineering Practice | proceeding | 9780784411902 | American Society of Civil | United States | Oceanography | Ocean Engineering | |
| 8905 | Coastal Hazards - Selected Papers from EMI 2010 | proceeding | 9780784412664 | American Society of Civil | United States | Safety, Risk, Reliability | | |
| 8908 | Coastal Sediments '07 - Proceedings of 6th International Symposium on Coastal Engineering and Science of Coastal Sediment Processes | proceeding | 9780784409268 | American Society of Civil | United States | Oceanography | Civil and Structural Eng | Ocean Engineering |
| 8909 | Coastal Structures 2003 - Proceedings of the Conference | proceeding | 9780784407332 | American Society of Civil | United States | Civil and Structural Eng | Ocean Engineering | Modeling and Simu |
| 8938 | Cold Regions Engineering Cold Regions Impacts on Transportation and Infrastructure: Proceedings of the Eleventh International Conferenc | proceeding | 9780784406212 | American Society of Civil | United States | Mechanical Engineerir | | |
| 9036 | Composite Construction in Steel and Concrete VI - Proceedings of the 2008 Composite Construction in Steel and Concrete Conference | proceeding | 9780784411421 | American Society of Civil | United States | Civil and Structural Eng | Mechanics of Material | Building and Cons |
| 9037 | Composite Construction in Steel and Concrete VII - Proceedings of the 2013 International Conference on Composite Construction in Steel a | proceeding | 9780784479735 | American Society of Civil | United States | Civil and Structural Eng | Mechanics of Material | Building and Cons |
| 9038 | Composites in Construction: A reality | proceeding | 9780784405963 | American Society of Civil | Italy | Civil and Structural Eng | Building and Construc | Ceramics and Com |
| 9052 | Computational Intelligence, From Theory to Practice - Proceedings of the 2008 Information Technology Symposium | proceeding | 9780784407608 | American Society of Civil | United States | Computer Science (all) | Civil and Structural Eng | Computational Ma |
| 9107 | Computing in Civil and Building Engineering | proceeding | 9780784405130 | American Society of Civil | United States | Computer Science (all) | Civil and Structural Eng | Building and Cons |
| 9108 | Computing in Civil and Building Engineering - Proceedings of the 2014 International Conference on Computing in Civil and Building Engine | proceeding | 9780784413616 | American Society of Civil | United States | Computer Science App | Civil and Structural Eng | Building and Cons |
| 9109 | Computing in Civil Engineering - Proceedings of the 2013 ASCE International Workshop on Computing in Civil Engineering | proceeding | 9780784477908 | American Society of Civil | United States | Civil and Structural Eng | Building and Construc | |
| 9113 | CONCREEP 2015: Mechanics and Physics of Creep, Shrinkage, and Durability of Concrete and Concrete Structures - Proceedings of the 10th | proceeding | 9780784479346 | American Society of Civil | United States | Mechanics of Materials | Building and Construc | |

DISCLAIMER-TERMS&CONDITIONS    SERIALS    CHINESE JPS on SERIALS LIST    NON-SERIALS    DISCONTINUED    DEFINITIONS

# Intelligent Hybrid Decision-Making for High-Speed Autonomous Driving Scenarios

Yifan Yang[1]; Yuchen Wu[2]; Gang Xu[3]; Yong Liu[4]; Zhitao Zhang[5]; and Jian Yang[6]

[1]Institute of Intelligent Transportation Systems, Polytechnic Institute, Zhejiang University, Hangzhou, Zhejiang, China. E-Mail: 22260198@zju.edu.cn
[2]Polytechnic Institute, Zhejiang University, Hangzhou, Zhejiang, China.
[3]College of Control Science and Engineering, Zhejiang University, Hangzhou, Zhejiang, China.
[4]Research Center for Intelligent Drive and Future Traffic, College of Control Science and Engineering, Zhejiang University, Hangzhou, Zhejiang, China (corresponding author). E-mail: yongliu@iipc.zju.edu.cn
[5]China Research and Development Academy of Machinery Equipment, Beijing, China.
[6]China Research and Development Academy of Machinery Equipment, Beijing, China.

## ABSTRACT

With the rapid development of autonomous driving technology, the transportation system is undergoing an unprecedented revolution. Due to the complexity of traffic rules and the real-time requirements of high-speed vehicles, decision-making techniques are of critical importance. This paper focuses on decision-making for high-speed autonomous driving, providing constraints for several common highway scenarios, including lane changes, overtaking, and navigating around accidents. In addition to making accurate decisions in these scenarios, autonomous vehicles must comply with traffic regulations and ensure smooth driving, avoiding sudden braking, sharp turns, and rapid acceleration. To address these challenges, this paper proposes a hybrid approach combining an improved deep reinforcement learning algorithm with rule-based control to design a decision-making algorithm for high-speed driving. Simulation results demonstrate that the proposed method improves decision accuracy by 25% compared to existing methods. These results highlight the strong applicability and generalization potential of the approach for real-world autonomous driving systems.

## INTRODUCTION

The rapid advancement of autonomous driving technology has positioned intelligent decision-making as a cornerstone for enhancing the safety and efficiency of transportation systems. Autonomous vehicles (AVs) must navigate increasingly complex traffic environments, characterized by high traffic density, dynamic road conditions, and high-speed scenarios. Effective decision-making frameworks are essential to address these challenges while ensuring compliance with traffic rules, smooth driving dynamics, and the safety of both passengers and surrounding road users.

In recent years, diverse decision-making models have emerged to tackle specific challenges in autonomous driving. For instance, reinforcement learning-based approaches have demonstrated significant potential in managing urban driving scenarios by optimizing tasks like lane changes and overtaking (Chen et al., 2018). Similarly, hybrid frameworks integrating deep reinforcement learning with model predictive control have shown promise in handling highway driving complexities, enhancing both decision accuracy and driving smoothness (Zhang et al., 2021).

Furthermore, addressing the interactions between autonomous and human-driven vehicles remains a pivotal challenge. Studies emphasize the importance of AVs predicting and adapting to the intentions of human drivers to ensure safety and seamless integration into mixed-traffic environments. For example, AVs are typically more conservative at higher speeds on arterials compared to human-driven vehicles, which influences their interaction dynamics with manual vehicles (Sinha et al., 2021). Another study explored how human drivers experience autonomous vehicles, noting that experienced drivers often prefer conventional vehicles, while novice drivers are more likely to trust AVs for their ease and safety (Manawadu et al., 2015).

Moreover, collaborative multi-agent systems and advanced communication protocols are essential to ensure safe interactions between autonomous and human-driven vehicles, particularly in mixed-traffic environments. Research has highlighted how such systems can leverage road infrastructure, vehicle-to-vehicle communications, and online motion prediction to enhance the safety and efficiency of autonomous driving (Aoki et al., 2021).

Despite significant progress, challenges persist, particularly in integrating holistic decision-making across diverse scenarios, managing the unpredictability of human-driven vehicles, and enabling robust communication in multi-agent frameworks. This paper addresses these gaps by proposing a hybrid decision-making strategy that combines rule-based mechanisms with advanced reinforcement learning to achieve adaptive, safe, and efficient decision-making for high-speed autonomous driving scenarios. Our main contributions are as follows:

• **Hybrid Approach**: Integration of rule-based decision-making with advanced deep reinforcement learning techniques to address the complexities of real-time decision-making in diverse and unpredictable traffic scenarios.

• **Guided Reward Reinforcement Learning Algorithm**: Development of a reinforcement learning algorithm with guided rewards to encourage safe, rule-compliant, and efficient driving behavior.

• **Comprehensive Testing**: Implementation of diverse highway scenarios, including real-world accident scenarios, in a simulation environment to thoroughly evaluate decision-making abilities.

## METHODOLOGY

### Scenarios

This study evaluates the decision-making capabilities of autonomous vehicles in multiple complex traffic scenarios, each designed to simulate real-world highway

driving conditions. The ego vehicle is required to avoid collisions, stay within lane boundaries, comply with speed limits, and prevent severe discomfort (e.g., sudden braking and sharp turns). The key scenarios include:
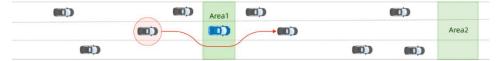
**Straight road yielding:** The ego vehicle travels from Area 1 to Area 2, encountering other vehicles that may attempt to overtake from the rear or side at random times, speeds, and positions. The vehicle must ensure safe and efficient arrival while avoiding collisions, as illustrated in Figure 1(a).

**Accident scene detouring:** The ego vehicle encounters an accident (or a broken-down vehicle) blocking the road, requiring a reasonable detour. The surrounding traffic density, speed, and detour direction are randomized, as shown in Figure 1(b).

**Low-speed vehicle overtaking:** The ego vehicle encounters a slow-moving vehicle (traveling at less than 30% of the speed limit) and must safely and efficiently overtake while considering the behavior of other road users, as depicted in Figure 1(c).

**Same-lane construction scenario:** The ego vehicle must navigate around a construction site demarcated by cones while driving from Area 1 to Area 2 under varying traffic densities, as shown in Figure 1(d).

These scenarios are intended to evaluate the robustness and adaptability of the proposed decision-making strategy in handling common high-speed driving situations, including overtaking, accident avoidance, and construction zones.



(a) Yield to other overtaking vehicles and arrive safely and efficiently



(b) Bypass the accident site and arrive safely and efficiently



(c) Overtake the low-speed vehicle ahead and arrive safely



(d) Bypass the same-lane construction area and arrive safely

**Figure 1. Key scenarios in high-speed autonomous driving**

**Research objectives and method selection**

The primary objective of this study is to develop a reliable decision-making strategy for high-speed autonomous driving that ensures accuracy, compliance with traffic regulations, and smooth driving dynamics. Compared to traditional decision-making algorithms, the Proximal Policy Optimization (PPO) reinforcement learning algorithm can effectively coordinate the relationships between various rules and decisions by adjusting factors such as reward functions, thereby achieving optimal results. One major advantage of the PPO algorithm is its simplicity in implementation, while achieving results comparable to or even better than other state-of-the-art algorithms, such as A2C, TRPO, and ACER, especially in navigation problems.

Deep reinforcement learning (DRL) is used to train models that adaptively learn optimal driving behaviors by interacting with the environment. Specifically, Categorical PPO is used to control the vehicle's steering, speed, and acceleration. These DRL methods effectively optimize the decision-making strategy through trial and error, maximizing safety and efficiency.

Rule-based control complements DRL by providing deterministic and interpretable control mechanisms, especially to mitigate the shortcomings of purely learning-based approaches. For example, PID control is employed to stabilize lane following after lane changes and to reduce jerks and large steering angles.

By combining these two methods, the system aims to leverage the adaptability of DRL while ensuring reliability and stability through traditional control mechanisms.

**Reinforcement learning framework**

The PPO agent interacts with the simulated environment and learns to make optimal decisions through repeated trials, aiming to maximize a cumulative reward that represents safe, efficient, and smooth driving behavior. To further enhance learning efficiency, a dynamic reward function is used, incorporating elements like collision avoidance, speed compliance, and passenger comfort. This reward function guides the agent towards desirable behaviors and penalizes unsafe actions, such as sudden lane changes or excessive acceleration. The framework diagram of the deep reinforcement learning model training phase is shown in Figure 2.



**Figure 2. Reinforcement learning training flowchart**

**Action space:** Training with a discrete action space is more straightforward and feasible, particularly in the context of autonomous driving, where the primary value of reinforcement learning lies in its ability to generate high-level decisions rather than perform low-level controls. High-level control focuses on tasks such as selecting lane changes or speed adjustments, rather than directly outputting continuous variables like acceleration or steering angles. Given the inherent lack of interpretability in neural network-based decision-making, deploying reinforcement learning algorithms for real-world applications becomes significantly more practical at the high-level decision-making layer. This facilitates seamless integration with other control algorithms. In contrast, directly applying reinforcement learning to low-level control tasks often encounters challenges in deployment feasibility.

To expedite convergence during training, a discrete action space is adopted. In a simplified two-dimensional simulation environment, the discrete control actions for the vehicle include: accelerate, decelerate, maintain speed, change to the left lane, change to the right lane, adjust acceleration, and modify steering angle, shown in the following parameters.

$$a_{space} = \left\{ s_{acc}, s_{dec}, v_t, d_{t \to left}, d_{t \to right}, a_t, r_t \right\}$$

**State space:** Considering that a single-frame observation cannot adequately represent the high-level driving behavior of surrounding vehicles, the state at time *t* is represented using the observational data from the previous five frames, up to time *t*. This approach effectively captures temporal dependencies and provides a more comprehensive representation of the driving environment. The details of the state space are presented in Table 1.

**Table 1. Main vehicle and obstacle status information**

| Vehicle Information | Obstacle Information |
|---|---|
| • Deviation between the target area position and the vehicle's current position | • Deviation between the obstacle position and the vehicle's current position |
| • Vehicle's heading angle | • Obstacle's heading angle |
| • Longitudinal velocity of the vehicle's rear axle | • Obstacle's velocity |
| • Longitudinal acceleration of the vehicle's rear axle | • Obstacle's acceleration |
| • Lateral acceleration of the vehicle's rear axle | • Obstacle's width |
| • Steering angle of the front wheels | • Obstacle's length |
| • Steering angle of the front wheels in the previous step | • None |
| • Longitudinal acceleration of the vehicle in the previous step | • None |
| • Current lane offset of the vehicle | • None |

**Network design:** The overall network design is illustrated in Figure 3, where the orange sections represent the vehicle's own information, and the yellow sections correspond to information about surrounding obstacles. Two separate normalization (Norm) modules are employed to standardize the vehicle's own data and obstacle data independently. The processed information is then passed through a feature extraction network, structured as a multi-layer perceptron (MLP), to extract features. The output features are combined to form a vector that represents the feature of $s_t$ which is subsequently used as input by both the Actor and Critic networks.



**Figure 3. The Overall network design**

The detailed structure of the Feature Net2, depicted in Figure 4, processes obstacle observations individually, regardless of the time step. Each obstacle observation is mapped independently through an embedding layer. The mapped results are then concatenated based on their respective time steps, preserving the original order. This concatenated output is further passed through another embedding layer to generate five vectors, each representing the aggregated obstacle features at a specific time step. These vectors are then concatenated again, and the resulting feature is mapped once more to obtain a representation of obstacle features over the past five frames. This temporal representation captures the dynamic environment effectively, enabling more informed decision-making.



**Figure 4. Feature Net2 design**

**Reword design:** to enable the vehicle to perform complex actions across various scenarios and reach the desired target point, a composite reward function is designed. The total reward is composed of three components: a basic reward, a collision avoidance reward, and a rule compliance reward. These components are combined to form the total reward, as shown in the following equation.

$$r_{total} = r_{base} + \alpha \times r_{collide} + \beta \times r_{rule}$$

The reward design includes the following considerations: (1) No collision penalty is applied during terminal states in training. (2) The longitudinal collision avoidance reward only considers obstacles/vehicles in front of the ego vehicle. (3) The coefficients $\alpha$ and $\beta$ correspond to the collision avoidance reward and rule compliance reward, respectively. As the training progresses, the weight of the rule compliance reward is gradually increased, guiding the autonomous vehicle to learn more complex behaviors. This dynamic reward adjustment aims to transition the vehicle from merely reaching the destination to doing so while adhering to traffic rules. In the early stages of training, the vehicle primarily relies on the basic reward to learn effective obstacle avoidance and to successfully reach the target. During the later stages, the rule compliance reward weight is dynamically adjusted, encouraging the vehicle to achieve stable and compliant driving behaviors. The detailed structure of the reward design is illustrated in Table 2.

**Table 2. Composite reward design**

| Reward | Items | Value |
|---|---|---|
| **Basic rewards** $(r_{base})$ | Step reward | -0.1 |
| | Distance reward | $0.5 \times \Delta d$ |
| | Arrival rewards | 200 |
| | Collision rewards | -200 |
| | Overtime reward | -200 |
| **Collision avoidance rewards** $(r_{collide})$ | Longitudinal collision avoidance | $\min(-0.1 \times (1 + (d_{ysafe} - d_x)), 0)$ |
| | Lateral collision avoidance | $\min(-0.1 \times (1 + (d_{xsafe} - d_y)), 0)$ |
| **Rule rewards** $(r_{rule})$ | Overspeed | $-0.1 \times \max(v - v_{rule}, 0)$ |
| | Large steering | $-\max(-0.4 - (|jerk_y| - 0.7) \times 3, -1.0)$ |
| | Sudden braking/acceleration | $-\max(-0.4 - (|jerk_y| - 0.7) \times 3, -0.9)$ |
| | Keep Lane centerline | $-\max(offset - 0.5, 0)$ |

**PID control and rules for enhancing model performance**

Although Categorical PPO has achieved commendable success rates, trajectory playback and data analysis revealed issues arising from the use of a discrete action space, which introduces coarse action granularity. This granularity results in instability, including noticeable jitter and difficulty maintaining straight-line driving on straight roads. Additionally, significant overshooting occurs during directional adjustments, leading to sharp steering maneuvers. To address these challenges, a PID controller was incorporated to enhance stability during straight-line driving.

Furthermore, action masking mechanisms were introduced to improve safety by preventing unintended or unsafe actions. The detailed implementation is illustrated in Figure 5.

Reinforcement learning is responsible for controlling the vehicle's steering and speed, while the PD controller is tasked with ensuring stability after lane changes. To achieve seamless integration between reinforcement learning and the PD-based control scheme, a hybrid strategy was designed to maintain stability across various lane control scenarios. The integration strategy is outlined as follows:

**Turning → Straight Driving:** When the current lane offset reaches the predefined threshold for alignment, the PD controller takes over the steering angle to stabilize the vehicle.

**Straight Driving → Turning:** When the current distance to the front vehicle falls below the specified safety threshold, reinforcement learning assumes control to manage obstacle avoidance and lane changes.

**Straight Driving → Straight Driving (Speed Suppression):** When the safe time to collision is less than the defined critical threshold, rule-based control overrides acceleration to execute emergency braking.

This hybrid approach ensures the stability and safety of the vehicle across various driving scenarios while leveraging the strengths of both reinforcement learning and rule-based controls.



**Figure 5. Add PD controller for model enhancement**

**SIMULATIONS AND RESULTS**

In autonomous driving simulation experiments, our proposed approach was tested in both our integrated simulation environment and the existing Highway-env environment. As illustrated in Figures 6(a), 6(b), 6(c), and 6(d), the ego vehicle demonstrated the ability to make accurate decisions in various randomized traffic scenarios. These scenarios included straight-road yielding, detouring around low-speed vehicles, and navigating same-lane construction zones. The results highlight the effectiveness of our hybrid decision-making framework in handling complex traffic conditions and making safe, efficient, and context-aware driving decisions. This hybrid approach ensures the stability and safety of the vehicle across various driving scenarios while leveraging the strengths of both reinforcement learning and rule-based controls.

The experimental results demonstrate the advantages of the proposed reinforcement learning framework with compound rewards over the standard PPO algorithm. Specifically, six experiments were conducted to evaluate the performance in terms of average driving time, success rates, and lane stability under various highway scenarios.

Figure 7(a) illustrates the comparison of average driving time between the standard PPO and PPO with compound rewards. The compound reward PPO consistently achieved shorter driving times across all experiments, with an average improvement of approximately 30%. This indicates the effectiveness of compound rewards in optimizing high-speed decision-making, enabling vehicles to navigate efficiently while maintaining safety and stability.



(a) Overtaking in integrated environment

(b) Bypass obstacle areas

(c) Avoid other vehicles in Highway-env

(d) Overtaking in Highway-env

**Figure 6. Simulation result**

In addition, the success rate comparison (Figure 7(b)) highlights a substantial improvement with the enhanced model. The success rate increased from an average of 61% with standard PPO to 86% with the incorporation of compound rewards, PID control, and rule-based adjustments. This improvement underscores the stability and robustness of the hybrid approach in handling complex traffic scenarios, such as lane changes and obstacle avoidance.

Figure 7(c) compares the steering angle variations between the two methods. The hybrid approach significantly reduced sharp steering adjustments, minimizing the risk of unstable maneuvers. The integration of PID control and action masking effectively addressed the limitations of coarse action granularity in discrete action spaces, ensuring smoother and safer driving dynamics.

Furthermore, Figure 7(d) highlights the improvement in maintaining lane stability during straight-line driving. The standard PPO exhibited frequent oscillations when controlling the vehicle on straight roads, resulting in deviations from the lane

center. In contrast, the proposed hybrid approach, with the integration of PID control for straight-line stability, effectively eliminated oscillations, ensuring smooth and jitter-free motion. This improvement demonstrates the ability of the framework to maintain precise control over the vehicle's trajectory, even in high-speed environments.

Overall, these results demonstrate that the proposed hybrid framework, combining reinforcement learning, PID control, and rule-based mechanisms, not only enhances the performance of autonomous driving systems in high-speed environments but also ensures a balance between efficiency, safety, and driving stability.



(a) Comparison of average driving time



(b) Comparison of success rate



(c) Comparison of steering angle



(d) Comparison of lane stability

**Figure 7. Simulation result**

**CONCLUSION**

This study proposed a novel hybrid decision-making framework for high-speed autonomous driving, integrating reinforcement learning, PID control, and rule-based mechanisms to effectively tackle the challenges of complex traffic scenarios. By testing the framework in both custom-built and standardized simulation environments, we demonstrated significant improvements in decision accuracy, driving stability, and compliance with traffic regulations. Specifically, the integration of PID control reduced oscillations during straight-line driving, ensuring smoother and more stable vehicle motion. Rule-based adjustments played a critical role in enhancing safety by mitigating abrupt maneuvers, while the compound reward structure allowed reinforcement learning to achieve higher success rates and faster driving times, striking an optimal balance between efficiency and safety. Comprehensive testing in challenging scenarios such as straight-road yielding, low-speed vehicle detouring, and construction zone navigation confirmed the

framework's robustness, effectiveness, and generalization capabilities across various driving conditions. These findings underscore the significant potential of combining data-driven and rule-based methods to create more reliable, efficient, and safe autonomous driving solutions. The contributions of this study provide a strong foundation for future work, which will extend the framework to more dynamic, cooperative, and real-world driving environments.

## REFERENCES

Chen, J., Tang, C., Xin, L., Li, S., & Tomizuka, M. (2018). Continuous decision making for on-road autonomous driving under uncertain and interactive environments. 2018 IEEE Intelligent Vehicles Symposium (IV), 1651-1658.

Zhang, W., Liu, Q., & Sun, Z. (2021). Hybrid decision-making framework for highway autonomous driving using deep reinforcement learning and model predictive control. IEEE Transactions on Vehicular Technology, 70(3), 2498-2511.

Sinha, A., Radwan, A., & Dixit, V. (2021). Interactions between human-driven and autonomous vehicles on public roads. arXiv.

Manawadu, U. E., Ishikawa, M., Kamezaki, M., & Sugano, S. (2015). Analysis of individual driving experience in autonomous and human-driven vehicles using a driving simulator. 2015 IEEE International Conference on Advanced Intelligent Mechatronics (AIM), 299-304. Manawadu et al., 2015.

Aoki, S., Lin, C.-W., & Rajkumar, R. (2021). Human-robot cooperation for autonomous vehicles and human drivers: Challenges and solutions. IEEE Communications Magazine, 59(1), 35-41.

Lefevre, S., Vasquez, D., & Laugier, C. (2016). "A survey on motion prediction and risk assessment for intelligent vehicles." Robomechanics and Automation Journal, 28(3), 3-19.

Kuutti, S., Fallah, S., Bowden, R., & Barber, P. (2020). "A survey of deep learning applications to autonomous vehicle control." IEEE Transactions on Intelligent Transportation Systems, 21(2), 712-733.

Hubmann, C., Becker, M., Althoff, D., Lenz, D., & Stiller, C. (2017). Decision Making for Autonomous Driving Considering Interaction and Uncertain Prediction of Surrounding Vehicles. 2017 IEEE Intelligent Vehicles Symposium (IV), 1671-1678.

Cao, Y., Chen, Y., & Liu, L. (2022). Research Prospect of Autonomous Driving Decision Technology under Complex Traffic Scenarios. MATEC Web of Conferences, 350, 03031.

Yildirim, M., Mozaffari, S., McCutcheon, L., Dianati, M., & Tamaddoni-Nezhad Saber Fallah, A. (2022). Prediction Based Decision Making for Autonomous Highway Driving. 2022 IEEE 25th International Conference on Intelligent Transportation Systems (ITSC), 138-145.

Lu, X., Zhao, H., Gao, B., Chen, W., & Chen, H. (2022). Decision-Making Method of Autonomous Vehicles in Urban Environments Considering Traffic Laws. IEEE Transactions on Intelligent Transportation Systems, 23, 21641-21652.

(54) 发明名称
　　一种多智能体高效协同路径规划方法

(57) 摘要

本发明涉及多智能体协同规划技术领域,具体公开了一种多智能体高效协同路径规划方法,包括系统初始化步骤、候选点集合扩充步骤、AStar算法路径规划步骤、任务执行监控步骤、候选点列表更新步骤、MILP目标点选择步骤、过渡目标点选择步骤和路径规划完成步骤,本发明能够综合考虑任务区域、智能体的初始位置、速度限制、目标位置坐标、障碍区域以及最少执行次数指标等因素,规划出各智能体的最优路径。通过合理的任务分配和障碍规避策略,确保智能体能够高效、安全地完成既定任务,在路径长度、任务完成时间和计算时间等关键性能指标上均优于传统方法。

**Algorithm 1** 多智能体高效协同路径规划方法

**Input:** 任务区域,智能体位置、速度和可执行次数信息,目标信息
**Output:** 每个智能体的路径和执行任务情况

1: 对于每个目标增加过渡点
2: **while** 未完成全部目标任务 **do**
3: 　利用AStar算法计算每个智能体到每个目标的路径和时间
4: 　给出每个智能体的满足时间和距离约束的目标
5: 　**if** 存在智能体没有满足约束的目标 **then**
6: 　　考虑过渡点
7: 　利用MILP规划所有智能体到目标路径长度之和最小的下一步目标
8: 　某个智能体最先完成查证
9: 　**if** 该智能体执行次数为正且该目标的任务指标已完成 **then**
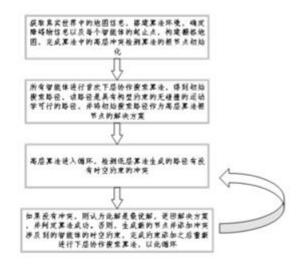10: 　　执行该目标对应任务
11: **return** 每个智能体的路径点和执行任务情况

（54）发明名称
　　一种基于分层搜索的多无人车构型保持协作运动规划方法

（57）摘要

本发明公开了一种基于分层搜索的多无人车构型保持协作运动规划方法，本发明通过分层搜索算法，在不涉及到具体编队控制的前提下实现了构型约束下的多无人车的运动规划，生成具有构型约束的无碰撞的运动学可行的路径，此路径能够在真实世界中有非常好的应用，同时构型形状、构型无人车数量可以任意指定，且可以实现构型变换、构型保持、构型分散、构型单车混合规划等多种功能，本发明搜索过程中加入了朝向角相同约束，具有朝向角尽量保持一致的特性，在探索的过程中可以尽可能让队里的无人车的朝向保持一致，防止出现无人车掉队、朝向偏差过大难以维持构型、搜索过程中信息检索不完全等情况，算法普适性强。

CN 118466486 A