

# 基于边缘云计算的 AI 一体机工具设计与实现

何俊

上海市信息网络有限公司

DOI:10.12238/acair.v3i1.11900

**[摘要]** 随着物联网技术和人工智能的快速发展,边缘计算与云计算的协同工作模式逐渐成为智能系统设计的新趋势。鉴于此,本文探讨了边缘计算与云计算协同驱动下的AI一体机工具的设计与实现,分析了其理论基础与技术背景,并详细介绍了系统设计、实现步骤以及性能评估方法。通过上海教师教育学院推训边缘云项目的应用案例,展示了AI一体机在提供高效能边缘计算、智能化教学资源管理以及构建稳定虚拟化平台方面的能力。

**[关键词]** 边缘计算; 云计算; AI一体机; 工具设计

中图分类号: TH-39 文献标识码: A

## Design and Implementation of an AI All-in-One Device Based on Edge Cloud Computing

Jun He

Shanghai Information Network Co., Ltd.

**[Abstract]** With the rapid development of Internet of Things (IoT) technology and artificial intelligence (AI), the collaborative working mode of edge computing and cloud computing has gradually become a new trend in the design of intelligent systems. In view of this, this paper discusses the design and implementation of an AI all-in-one machine tool driven by the synergy between edge computing and cloud computing, analyzes its theoretical foundation and technical background, and introduces in detail the system design, implementation steps, and performance evaluation methods. Through the application case of the Edge Cloud project promoted by Shanghai Academy of Educational Sciences for teachers' training, it demonstrates the capabilities of the AI all-in-one machine in providing high-efficiency edge computing, intelligent management of educational resources, and building a stable virtualization platform.

**[Key words]** Edge Computing; Cloud Computing; AI All-in-One Machine; Tool Design.

### 引言

在数字化时代,智能技术快速发展,AI应用需求不断增长。边缘计算与云计算协同模式逐渐兴起,为解决数据处理效率、隐私安全等问题提供新思路<sup>[1]</sup>。AI一体机作为该协同模式下的产物,集成多种先进技术,其硬件架构、软件平台及协同机制的合理设计成为关键。在教育等领域,AI一体机有望成为推动创新与发展的重要力量,对其深入研究与应用探索极具价值。

### 1 理论基础与技术背景

边缘计算通过在数据源附近进行数据处理,能够显著降低延迟,提高效率,这对于需要实时响应的应用场景如自动驾驶至关重要。同时,边缘计算有助于节省云端带宽资源,减轻中心服务器压力,实现绿色计算<sup>[2]</sup>。在隐私保护方面,边缘计算能够在本地处理敏感数据,增强用户信任。系统可靠性也因边缘计算而增强,即使在互联网连接不稳定的情况下,关键服务也能得到保障。技术实现上,多样化的硬件平台如Raspberry Pi和NVIDIA

Jetson系列为Edge AI提供了物理支撑。软件框架如TensorFlow Lite和OpenVINO使得模型能在资源受限的边缘设备上高效运行。在边缘智能芯片方面,基于轻量级CPU、GPU、ASIC和FPGA的硬件为深度学习应用提供了专门设计的处理器<sup>[3]</sup>。

### 2 AI一体机系统设计

#### 2.1 硬件架构与优化

AI一体机的硬件架构是其性能和效率的基础,旨在为边缘计算和云计算的协同工作提供强有力的支持。首先,该架构集成了高性能处理器(如GPU或TPU),这些处理器专为加速深度学习算法而设计,能够快速处理复杂的数学运算,从而显著提升模型推理速度<sup>[4]</sup>。此外,为了适应不同应用场景的需求,硬件还配备了丰富的接口和扩展槽,支持多种传感器和其他外设的连接,确保数据采集的多样性和灵活性。更重要的是,考虑到功耗和散热问题,硬件设计中融入了节能技术和高效的散热方案,使得AI一体机可以在各种环境下稳定运行。同时,硬件层面的安全防护措

施也不可或缺,通过内置的安全芯片和加密模块,保障数据传输的安全性,防止未经授权的访问或篡改。总之,精心设计的硬件架构不仅提升了AI一体机的整体性能,也为后续软件开发提供了坚实的平台支撑。

### 2.2 软件平台与工具链

在AI一体机系统设计中,软件平台与工具链扮演着至关重要的角色,它们决定了系统的易用性和开发效率。一个完善的软件平台应具备强大的开发环境、丰富的API接口以及直观的用户界面。具体来说,平台需集成最新的深度学习框架(如TensorFlow、PyTorch等),以便开发者可以轻松构建和训练复杂的神经网络模型<sup>[5]</sup>。与此同时,针对边缘计算的特点,平台还需提供轻量级的推理引擎,以确保模型能够在资源受限的环境中高效运行。

### 2.3 云边协同机制

AI一体机的成功运作离不开高效的云边协同机制,这种机制确保了边缘侧和云端之间的紧密合作,最大化地发挥了各自的优势。首先,通过采用先进的网络通信协议(如MQTT、CoAP等),边缘设备可以实时与云端交换数据,实现低延迟、高可靠性的信息传递。当本地处理能力不足或需要更大规模的数据分析时,边缘设备会将任务卸载至云端,利用后者强大的计算资源完成复杂计算。反之,在云端训练好的模型也可以推送到边缘端进行推理,这样既减少了带宽占用,又提高了响应速度<sup>[6]</sup>。此外,云边协同还包括智能任务调度策略,根据实际需求动态分配计算资源,避免资源浪费。这样的协作模式不仅提升了整体系统的灵活性和适应性,也为用户带来了更好的体验。最后,为了保证系统的长期稳定运行,云边协同机制还需要考虑数据同步、故障恢复等功能,确保即使在网络中断的情况下也能维持基本服务,保障业务连续性。

## 3 系统实现与性能评估

### 3.1 系统实现步骤

首先进行硬件选型与组装,依据设计的硬件架构挑选适配的高性能处理器、存储设备及通信模块,搭建稳定的硬件平台,并完成操作系统及底层驱动的安装与调试,确保硬件资源能被有效调用。接着进行软件框架搭建,整合各类AI算法库和开发工具,构建起分层、模块化的软件体系,实现各功能模块的代码编写与集成,重点攻克模块间的数据交互与接口适配问题,保证系统的完整性和连贯性<sup>[7]</sup>。随后进行系统联调,模拟多种实际场景下的数据输入和业务流程,对各模块协同工作情况进行全面测试与优化,及时发现并解决潜在的冲突和错误,确保系统稳定运行,最终完成AI一体机的系统实现,使其具备投入实际应用的条件。

### 3.2 性能指标选取

选取准确率、召回率等指标衡量AI模型的识别精度,以评估一体机在图像识别、语音识别等任务中的表现,确保其能精准地完成各类智能分析任务<sup>[8]</sup>。采用处理速度、吞吐量作为性能衡量标准,反映系统在单位时间内处理的数据量和完成任务的效

率,对于实时性要求高的场景至关重要。延迟指标用于评估系统响应的及时性,尤其是在边缘计算场景下,低延迟能保障系统快速反馈,提升用户体验。资源利用率方面,关注CPU、内存、GPU等硬件资源的使用情况,通过合理优化资源分配,在保证性能的同时避免资源浪费,使系统在不同负载下都能高效运行,这些指标从多个维度全面反映AI一体机的性能表现。

### 3.3 性能测试方法

采用模拟真实场景的数据集进行测试,如在图像识别中使用包含不同光照、角度、背景的图像集,确保测试数据的多样性和复杂性,能真实反映系统在实际应用中的性能表现。运用压力测试工具,逐步增加系统的负载,模拟高并发场景下的运行状况,监测系统在不同负载压力下的各项性能指标变化,确定系统的瓶颈和极限处理能力,评估其稳定性和可靠性<sup>[9]</sup>。同时,进行长时间的稳定性测试,让系统持续运行一段时间,观察其在长时间运行过程中的性能波动、资源泄漏等问题,以保证系统在实际生产环境中能够稳定、可靠地运行,为性能评估提供全面、准确的数据支持。

### 3.4 性能优化策略

针对性能瓶颈,优化硬件配置是一种直接有效的方法,例如升级GPU以加速计算密集型任务,增加内存容量以提升数据处理效率,确保硬件资源能够充分满足系统需求。在软件层面,通过算法优化减少计算复杂度,如采用轻量级的神经网络模型或优化算法结构,提高模型的推理速度<sup>[10]</sup>。同时,采用缓存机制,对频繁访问的数据进行缓存,减少数据读取时间,提高系统响应速度。优化数据传输流程,利用数据压缩、异步传输等技术,降低数据传输延迟和带宽占用,提升系统整体性能,使AI一体机在实际应用中能够达到更优的性能表现,满足不同场景下的业务需求。

## 4 项目应用案例分析

### 4.1 项目介绍

上海教师教育学院推训边缘云项目所用设备为两台AI推训一体机,利用一体机的推理和计算能力,搭建为教师教育学院专用的边缘云项目。

### 4.2 AI一体机在项目中的应用

#### 4.2.1 高效能的边缘计算支持教师教育学院的实训环境

在上海教师教育学院推训边缘云项目中,两台AI一体机凭借其卓越的硬件性能和软件集成能力,为教师教育学院提供了一个高度定制化的实训环境。每台AI一体机配备了高性能的GPU(如NVIDIA RTX 5880 ada),这些GPU拥有强大的并行处理能力,可以动态切割服务器上的GPU资源分配给不同的实训终端,实现服务器资源对终端的最大化算力资源支持。具体而言,单服务器(8卡GPU)最多可支持64个独立GPU实训终端,这意味着即使在一个班级中有大量学员同时进行深度学习或数据分析等高负载任务时,每位学员也能获得足够的计算资源。此外,AI一体机内置了丰富的课程素材管理和自动分发功能,系统会将最新的教学内容实时推送到不同账户的容器中,确保每位学生都能及

时获取最新资讯,无需长时间等待下载。这种高效的资源配置方式不仅提高了培训效率,还极大地改善了学员的学习体验,使得他们能够更加专注于实践操作而非被技术细节所困扰。

#### 4.2.2 智能的教学资源管理与辅助教学服务

针对传统教学资源管理中存在的私有分享、下载不易及版本管理难题,丽台AI一体机引入了一套智能化的教学资源管理系统。该系统通过AI助教问答式指导和大模型召回的方式,提供了教辅材料自助问答查询服务,允许用户上传文件,利用RAG(检索增强生成)技术进行知识库问答,从而更精准地解决老师的问题。同时,AI一体机还配备了丰富的AI教学备课材料库,包括详尽的课程讲义、生动的示例代码以及互动实验等,老师们无需再为编写教案而费心,只需一键调取即可获得结构清晰、内容新颖的专业教学素材。这样的设计不仅简化了教学准备过程,也为教师们提供了更多创新教学方法的可能性。更重要的是,这套系统还支持多轮对话功能,即对话包含上下文信息,这有助于维持连续性的讨论氛围,促进师生之间的深入交流,最终提升教学质量。

#### 4.2.3 构建稳定且易于维护的虚拟化平台

为了满足教师教育学院对于稳定性、安全性和易用性的严格要求,丽台AI一体机构建了一个稳定且易于维护的虚拟化平台。这个平台基于Linux生态(如Ubuntu/信创系统),并集成了AWS Cloud Kubernetes Docker等先进技术,实现了从操作系统层到应用中间层再到业务应用层的全面覆盖。通过MaaS微调服务、推理服务等功能模块,AI一体机能够快速部署和持续升级各种AI应用,确保系统始终保持最新状态。特别是在数据安全方面,AI一体机采用了数据私有化定时备份机制,定期将重要数据备份到云端或本地存储设备,防止意外丢失。此外,考虑到未来可能的需求变化和技术进步,AI一体机预留了充足的扩展空间,支持全国产硬件和国产操作系统,以适应不断发展的教育信息化趋势。

#### 4.3 经验总结

在上海教师教育学院推训边缘云项目中,丽台AI一体机的应用展示了其在教育领域的巨大潜力和显著成效。通过强大的硬件性能和智能软件集成,AI一体机不仅为学员提供了高效、稳定的实训环境,还大幅提升了教学资源管理和辅助服务的智能化水平。引入的AI助教问答系统和大模型召回技术有效解决了

传统教学资源管理中的难题。

## 5 结论

AI一体机作为边缘计算与云计算协同的重要产物,在推动AI应用落地方面发挥着重要作用。通过将高性能计算、AI算法和云服务整合到单一设备中,AI一体机能够为各行各业提供高效、便捷的AI解决方案。特别是在教育领域,丽台AI一体机的应用成功地改善了教师教育学院的教学实训环境,简化了教学资源管理流程,并构建了一个稳定且易于维护的虚拟化平台。这些成果不仅满足了当前教育信息化的需求,也为未来更广泛的应用场景提供了参考模板。随着边缘计算和云计算技术的不断发展,AI一体机的功能将更加完善,应用场景也将更加广泛,为人工智能技术的普及和发展做出更大的贡献。

## [参考文献]

- [1]潘志宇.基于云计算与边缘计算的电网智能运维平台研究[J].电工技术,2024,(S1):96-98.
- [2]李宁宁,金永均,苗梦奇,等.基于边缘智能与云计算的人工智能型智慧育苗系统研究[J].农业装备技术,2024,50(03):7-12.
- [3]朱斌,刘东,刘天元,等.边缘计算在电力系统供需互动应用的研究进展与展望[J].电网技术,2024,48(10):4327-4341.
- [4]周楠,蔡頔,刘凯伦,等.基于边缘云计算的移动端电力数据交互技术研究[J].电子设计工程,2024,32(10):87-91.
- [5]张森,张英威,席珺琳,等.智能物联网中的边缘计算与数据处理[J].通讯世界,2024,31(03):132-135.
- [6]李志勇,丁国强,于绍晨,等.算力时代融合边缘云建设思路探讨[J].电信工程技术与标准化,2024,37(01):8-11+30.
- [7]牛岚甲,吴迪,刘全明,等.边缘计算环境中基于区块链的物联网解决方案[J].南京理工大学学报,2023,47(05):678-684.
- [8]贺迎先.基于边缘云计算的混合并行AI训练机制研究[J].西安文理学院学报(自然科学版),2023,26(01):9-12+16.
- [9]周振.基于边缘云计算的数据智能云平台技术研究[J].信息与电脑(理论版),2021,33(09):25-27.
- [10]蒋林涛.云计算、边缘计算和算力网络[J].信息通信技术,2020,14(04):4-8.

## 作者简介:

何俊(1983-),男,汉族,上海人,研究生,中级工程师,研究方向:网络安全和算力网络。