

嵌入式AI算法在边缘计算设备中的实现与挑战

黄锦俊

广州中海达定位技术有限公司 广东广州 510000

摘要：随着边缘计算技术的快速发展，嵌入式AI算法凭借低延迟、高隐私性的优势，成为边缘设备智能化的核心支撑。本文围绕嵌入式AI算法在边缘计算设备中的实现与挑战展开研究，先分析嵌入式AI与边缘计算融合的技术基础，明确二者适配的核心条件；再梳理算法部署、模型优化、硬件适配的实现路径，细化各环节关键操作；接着探讨性能、能耗、兼容性方面的核心挑战，剖析问题产生的技术根源；最后提出针对性优化方向与未来发展思路。研究旨在为嵌入式AI算法在边缘设备中的高效落地提供参考，推动边缘智能技术在物联网、工业互联网等领域的实际应用。

关键词：嵌入式AI算法；边缘计算设备；算法实现；性能优化

引言

在物联网、工业互联网等领域的需求驱动下，边缘计算设备需具备实时数据处理与智能决策能力，以满足场景对响应速度与数据隐私的高要求。嵌入式AI算法能将AI功能集成到资源受限的边缘设备中，避免数据传输至云端带来的延迟与隐私泄露问题，成为边缘设备智能化的关键技术。当前，嵌入式AI算法在边缘设备中的应用逐渐广泛，但受限于边缘设备的算力有限、存储容量小、能耗约束严格等特点，算法实现过程中面临模型适配难、运行性能不足、跨设备兼容性差等问题。

一、嵌入式AI算法与边缘计算设备的融合基础

1. 嵌入式AI算法的核心技术特征与边缘计算适配性

嵌入式AI算法具备轻量化、低资源消耗、高实时性的核心技术特征，这与边缘计算设备的需求高度契合，形成良好适配性。轻量化特征使算法能在边缘设备有限的存储空间中部署，无需依赖大规模云端存储；低资源消耗特性适配边缘设备算力与能耗约束，可在不造成设备过载的前提下运行；高实时性则满足边缘计算对数据处理速度的要求，确保智能决策快速输出。这种适配性让嵌入式AI算法能有效弥补边缘设备“算力弱但需智能”的短板，成为边缘计算设备实现智能化的核心技术支撑，为二者融合奠定基础。

2. 边缘计算设备的硬件架构与资源约束条件

边缘计算设备的硬件架构多以低功耗处理器为核心，搭配小型化存储模块、专用传感器接口与通信模块，整

体呈现“精简高效”的特点。常见的硬件架构包括基于ARM架构的处理器设计、集成专用AI加速单元的芯片方案等，这些架构旨在平衡性能与能耗。但边缘设备存在明显的资源约束，算力方面，多数设备无法支撑复杂AI模型的大规模并行计算；存储方面，本地存储容量有限，难以容纳大体积模型文件与海量数据；能耗方面，设备多依赖电池供电或低功耗供电模式，对算法运行的能耗需求严格，这些约束条件直接影响嵌入式AI算法的实现方式。

3. 嵌入式AI与边缘计算融合的应用场景需求分析

嵌入式AI与边缘计算的融合在多个场景中存在明确需求，不同场景对技术的要求各有侧重。在工业物联网场景中，需求集中于设备故障实时检测，需嵌入式AI算法在边缘设备上快速分析传感器数据，及时识别故障信号；在智能家居场景中，需求聚焦于用户行为识别与场景联动，算法需在边缘设备上低延迟处理环境数据，实现设备智能控制；在智慧交通场景中，需求围绕车辆实时感知与决策，算法需在车载边缘设备上高效处理摄像头、雷达数据，保障行车安全。这些场景需求推动嵌入式AI与边缘计算深度融合，也为技术发展指明方向。

二、嵌入式AI算法在边缘计算设备中的实现路径

1. 嵌入式AI算法的轻量化模型设计与压缩方法

嵌入式AI算法的轻量化模型设计需从模型结构入手，通过简化网络层数、减少参数数量实现轻量化，例如采用深度可分离卷积替代传统卷积，降低模型复杂度。模型压缩方法则包括量化、剪枝、蒸馏等，量化将模型

参数从高精度转为低精度，减少存储与计算需求；剪枝去除模型中冗余的参数与神经元，保留核心计算单元；蒸馏通过训练小模型模仿大模型的输出，在保证精度的同时缩小模型体积。这些设计与压缩方法需结合边缘设备资源约束，在模型体积、计算量与精度之间找到平衡，确保算法能在边缘设备上高效运行。

2. 边缘设备端AI算法的部署流程与工具链应用

边缘设备端AI算法的部署流程包含模型转换、适配调试、验证优化三个核心环节。模型转换阶段，需将训练好的AI模型转换为边缘设备支持的格式，确保模型能被设备硬件识别；适配调试阶段，结合设备硬件特性调整模型运行参数，解决模型与硬件不兼容的问题；验证优化阶段，测试模型在设备上的运行性能，优化运行效率。工具链的应用贯穿整个流程，常见的工具包括Tensor Flow Lite、ONNX Runtime等，这些工具能提供模型转换、优化、部署的一站式支持，帮助开发者简化部署流程，降低技术门槛，确保算法顺利在边缘设备上落地。

3. 嵌入式AI与边缘设备硬件接口的适配技术实现

嵌入式AI与边缘设备硬件接口的适配需围绕数据传输、控制指令交互展开，确保算法与硬件高效协同。数据接口适配方面，需实现算法与设备传感器、存储模块的数据交互，通过标准化接口协议（如SPI、I2C）确保数据实时、准确传输；控制接口适配方面，需建立算法与设备处理器、AI加速单元的控制链路，让算法能调用硬件资源进行计算。适配技术实现中，需针对不同硬件架构优化接口驱动程序，解决数据传输延迟、资源调用冲突等问题，例如为集成AI加速单元的设备开发专用驱动，让嵌入式AI算法能高效利用加速单元提升计算性能。

三、嵌入式AI算法在边缘设备中的性能优化方向

1. 基于边缘设备算力的AI算法推理速度优化策略

基于边缘设备算力的AI算法推理速度优化，需从计算流程与资源调度两方面入手。计算流程优化方面，通过算子融合减少计算步骤，将多个独立算子合并为单一算子，降低数据读写次数；采用并行计算技术，利用边缘设备多核处理器的优势，将推理任务拆解为子任务并行处理。资源调度优化方面，优先为AI算法分配核心算力资源，避免其他任务占用关键硬件资源；动态调整算力分配策略，在算法推理高峰期提升算力供给，空闲期减少算力消耗。这些策略能充分利用边缘设备有限的算力，提升算法推理速度，满足实时性需求。

2. 面向边缘设备能耗约束的AI算法资源调度方法

面向边缘设备能耗约束的AI算法资源调度，需以“低能耗、高效能”为核心目标。在硬件资源调度上，根据算法运行阶段动态调整硬件工作状态，例如推理任务较轻时，让部分处理器核心进入低功耗模式；在软件资源调度上，优化算法运行流程，减少不必要的计算与数据传输，降低软件层面的能耗。同时，采用能耗感知的任务调度策略，将高能耗的AI计算任务安排在设备供电稳定的时段，低能耗任务在电池供电时段执行，平衡能耗与任务需求，确保算法在满足能耗约束的前提下正常运行。

3. 边缘场景下嵌入式AI算法的精度与效率平衡路径

边缘场景下嵌入式AI算法的精度与效率平衡，需结合场景需求动态调整优化方向。对于精度要求高的场景（如医疗诊断辅助），采用“精度优先、适度妥协效率”的策略，通过模型蒸馏、量化时保留关键层精度等方式，在大幅降低效率的同时保障精度；对于效率要求高的场景（如实时监控），采用“效率优先、精度合理让步”的策略，通过简化模型结构、合理剪枝等方式提升效率，同时确保精度满足场景基本需求。此外，可引入动态调整机制，让算法根据边缘场景数据特征与设备状态，实时切换精度与效率模式，实现二者动态平衡。

四、嵌入式AI算法在边缘计算设备中的核心挑战

1. 边缘设备硬件异构性导致的算法兼容性问题

边缘设备硬件架构呈现显著异构性，不同设备可能采用ARM、x86、RISC-V等不同架构的处理器，集成的AI加速单元也存在差异，部分设备甚至无专用加速单元。这种异构性导致嵌入式AI算法难以形成统一的适配方案，同一算法在不同硬件设备上可能出现运行异常、性能差异大等兼容性问题。例如，针对ARM架构优化的算法在x86架构设备上运行时，可能因指令集不兼容导致推理错误；依赖专用加速单元的算法在无加速单元的设备上无法发挥性能优势。硬件异构性增加了算法适配难度，也提高了开发与维护成本。

2. 嵌入式AI算法在多任务场景下的实时性保障难点

在多任务并发的边缘场景中，嵌入式AI算法需与设备上的其他任务（如数据采集、通信传输）共享硬件资源，实时性保障面临挑战。当多个任务同时竞争算力、存储资源时，AI算法的推理过程可能被打断或延迟，导致智能决策输出滞后，无法满足场景实时需求。例如，在车载边缘设备中，AI导航算法需与音乐播放、蓝牙通

信等任务共享资源，若资源分配不合理，导航算法的实时路径规划可能延迟，影响行车体验。此外，任务负载的动态变化也增加了实时性保障难度，难以提前制定固定的资源分配方案。

3. 边缘设备动态环境对AI算法稳定性的影响因素

边缘设备多部署在复杂动态环境中，环境因素会直接影响嵌入式AI算法的稳定性。温度变化方面，过高或过低的温度会导致设备硬件性能波动，影响算法计算精度与速度；电磁干扰方面，工业场景中的电磁辐射可能干扰设备数据传输，导致AI算法获取的输入数据失真，影响推理结果；网络波动方面，部分边缘设备需与云端或其他设备协同，网络延迟或中断会导致算法无法获取必要的辅助数据，影响运行稳定性。这些动态环境因素让嵌入式AI算法的运行状态难以保持稳定，增加了技术应用风险。

五、嵌入式AI算法在边缘计算设备中的发展对策与展望

1. 针对边缘设备特性的定制化AI模型开发思路

针对边缘设备特性的定制化AI模型开发，需建立“设备需求驱动”的开发模式。首先，深入分析目标边缘设备的硬件参数（算力、存储、能耗）与应用场景需求，明确模型设计的核心约束；其次，采用模块化设计方法，将模型拆分为适配不同硬件模块的子模块，提升模型与设备的适配性；最后，引入硬件感知的模型训练策略，在模型训练过程中融入设备硬件特性数据，让模型在训练阶段就适应设备运行环境。这种定制化开发思路能让AI模型更贴合边缘设备特性，充分发挥设备性能，提升算法运行效果。

2. 嵌入式AI与边缘计算协同优化的技术标准构建

构建嵌入式AI与边缘计算协同优化的技术标准，需从接口、性能评估、安全三个维度入手。接口标准方面，制定统一的模型格式、数据传输协议与硬件调用接口，解决算法与设备的兼容性问题；性能评估标准方面，建立涵盖算力利用率、能耗效率、推理延迟的评估指标体系，为技术优化提供统一依据；安全标准方面，制定边缘设备上AI算法的数据加密、权限管理规范，保障数据

与算法安全。技术标准的构建需联合行业企业、科研机构共同参与，确保标准的科学性与实用性，推动行业技术规范化发展。

3. 多领域边缘智能应用中嵌入式AI算法的创新方向

在多领域边缘智能应用中，嵌入式AI算法的创新方向需结合场景需求持续拓展。在工业领域，可向“多模态数据融合推理”方向创新，融合传感器、图像、声音数据提升故障检测精度；在消费电子领域，可向“低功耗轻量化模型”方向创新，满足设备长续航需求；在公共安全领域，可向“实时动态识别”方向创新，提升复杂场景下的目标识别与追踪能力。同时，随着边缘计算与5G、物联网技术的融合，嵌入式AI算法可向“跨设备协同推理”方向创新，通过多个边缘设备协同完成复杂AI任务，拓展技术应用边界。

结束语

嵌入式AI算法与边缘计算设备的融合是推动边缘智能技术发展的核心方向，二者的结合能有效满足物联网、工业互联网等领域对低延迟、高隐私性智能服务的需求。本文通过分析二者融合的技术基础，明确了嵌入式AI算法的适配性、边缘设备的资源约束与场景需求；梳理算法实现路径，涵盖模型轻量化、部署流程与硬件适配；探讨性能优化方向，从推理速度、能耗控制、精度效率平衡展开；剖析核心挑战，包括硬件异构性、多任务实时性与动态环境影响；提出发展对策，涉及定制化模型、技术标准与场景化创新。这些研究内容为嵌入式AI算法在边缘设备中的高效落地提供了系统思路。

参考文献

- [1] 简毅, 柳建, 徐灵飞, 等. 嵌入式边缘检测算法的HLS加速实现[J]. 电子设计工程, 2020, 28(22): 132-135+140.
- [2] 高彬. 边缘增强算法的Simulink-嵌入式一体化设计与实现[D]. 中国石油大学(华东), 2020.
- [3] 王光宁. 基于嵌入式GPU的模切工件边缘缺陷检测算法研发[D]. 浙江大学, 2020.