

摘要

大语言模型（LLM）在多个领域的成功应用推动了人工智能发展，但其存在计算成本高、内存需求大等缺点，优化其推理效率、降低资源消耗成核心问题。主流大模型多为 **Decoder-Only** 架构，其中矩阵向量乘算子（GEMV）在大模型推理中占主导地位。GEMV 因自身特点为内存瓶颈任务，GPU 等计算密集型硬件利用率低，在边端场景下常成系统瓶颈。近存计算是新兴计算范式，能解决传统计算架构内存瓶颈问题。新兴的近存计算硬件，有高带宽、高存储、高并行性、高能效等优势，也存在计算能力弱、通信开销大等局限。如何从软件角度适配硬件加速，以及如何修改设计硬件优化其应用性能的软硬协同优化方案成为待深入研究的课题。

本论文基于近存计算技术研究算子 GEMV 的软硬协同加速方案：1）在近存计算硬件上，设计基于查找表的矩阵向量乘算法，使用查找表消除复杂计算，对近存计算硬件中的多级存储结构分别做出了访存优化，包括查找表分块算法减少通过 DMA 访问内存的次数，以及矩阵的行列重排减少对缓存的访问，增强寄存器的数据局部性；2）在模拟器上，为解决真实硬件上软件算法的计算瓶颈，基于周期精确近存激素模拟器设计硬件及指令以增强其计算能力，包括融合查表和加法指令用于高效快速查询乘积并累加，以及设计基于查找表的向量指令进行向量化的查表和访存，提高计算效率。最后对上述算法进行了详细的测试，分别在 CPU 和 GPU 以及近存硬件三个硬件平台上，进行了包括对 GEMV 算子总吞吐和能效比的测试。实验结果表明，在近存计算硬件上 GEMV 算子吞吐为 355GOPS，大概达到理论性能的 93%，是 CPU 平台的 14.7 倍，能效比大概是 CPU 平台的 8.6 倍，是 GPU 平台的 1.13 倍。

关键词：近存计算；矩阵向量乘；查找表

Abstract

The successful application of Large Language Models (LLMs) in multiple domains has propelled the advancement of artificial intelligence. Nevertheless, they are plagued by drawbacks such as high computational costs and substantial memory demands. Optimizing the inference efficiency and reducing resource consumption have emerged as core issues. Mainstream LLMs adopt the Decoder-Only architecture, within which the Generalized Matrix-Vector Multiplication (GEMV) holds a dominant position in LLM inferences. Due to its inherent characteristics, GEMV constitutes a memory bounded task, with the utilization of computing-intensive hardware like GPUs being low. It has always become a system bottleneck in edge scenarios. Processing-in-Memory (PIM) represents an emerging computing paradigm capable of resolving the memory bounded problem in traditional computing architectures. Emerging PIM hardware boasts advantages such as high bandwidth, high storage capacity, high parallelism, and high energy efficiency, yet it also suffers from limitations including weak computing capabilities and significant communication overhead. How to adapt hardware acceleration from a software perspective and how to modify the design of hardware to optimize its application performance through a hardware-software co-design solution have become subjects warranting in-depth research.

This thesis investigates the hardware-software co-design solution for GEMV based on Processing-in-Memory technology: 1) On the PIM hardware, a matrix-vector multiplication algorithm based on lookup tables is devised. Lookup tables are employed to eliminate complex computations, and memory access optimizations are implemented for the multi-level storage structure within the PIM hardware. This encompasses the lookup table blocking algorithm to reduce the number of memory accesses via DMA and the rearrangement of matrix rows and columns to minimize cache access and en-

hance data locality in registers. 2) On the simulator, to address the computational bottlenecks of software algorithms on real hardware, special processing unit and instructions are designed based on the cycle-accurate PIM simulator to enhance its computing capacity. This includes the fusion of table lookup and addition instructions for efficient and rapid query of products and accumulation, as well as the design of vector instructions based on lookup tables for vectorized lookup and memory access to boost computing efficiency. Finally, detailed tests were carried out on the aforementioned algorithms on three hardware platforms, namely CPU, GPU, and PIM hardware, including tests on the total throughput and energy efficiency of the GEMV operator. The experimental results indicate that the throughput of the GEMV operator on the PIM hardware amounts to 355 GOPS, approximately reaching 93% of the theoretical performance. It is 14.7 times that of the CPU platform, and the energy efficiency is approximately 8.6 times that of the CPU platform and 1.13 times that of the GPU platform.

Key Words : Processing-in-Memory; GEMV; Lookup-Table

目录

第 1 章 绪论	1
1.1 研究背景和意义	1
1.2 国内外研究现状	2
1.3 研究内容和创新点	4
1.4 论文组织结构	7
第 2 章 相关工作	8
2.1 大模型推理介绍	8
2.1.1 基本概念	8
2.1.2 大模型量化技术	11
2.1.3 矩阵向量乘算子	13
2.2 近存计算研究现状	15
2.2.1 近存计算技术发展	15
2.2.2 近存硬件 UPMEM	17
2.2.3 UPMEM 相关应用	20
2.3 本章小结	21
第 3 章 基于近存计算硬件的矩阵向量乘软件优化	22
3.1 基于多级存储的查找表分块算法	22
3.1.1 矩阵向量乘查找表基础	23
3.1.2 分块载入查找表卸载乘法	24

3.1.3 基于数制映射表卸载加法	26
3.1.4 算法小结	28
3.2 基于缓存的矩阵行列重排算法	29
3.2.1 矩阵行重排.	29
3.2.2 矩阵列重排.	30
3.3 本章小结	34
第 4 章 近存计算架构中的矩阵向量乘硬件设计	35
4.1 软件优化性能瓶颈分析	35
4.2 查表和加法融合的指令设计.	36
4.3 基于查表的向量指令设计实现	38
4.4 基于近存模拟器的硬件设计实现	42
4.5 本章小结	44
第 5 章 实验结果与分析	45
5.1 环境配置介绍.	45
5.1.1 硬件平台	45
5.1.2 数据准备和矩阵尺寸选择	46
5.1.3 基线设置	46
5.2 基于近存计算平台的软件优化测试与分析	47
5.2.1 总吞吐对比测试	47
5.2.2 能效比测试与瓶颈分析.	49
5.2.3 扩展性测试.	51
5.3 基于模拟器平台的硬件设计测试与分析.	53
5.4 本章小结	56
第 6 章 总结与展望	57
6.1 总结.	57

6.2 展望	58
参考文献	59
致谢	70

图目录

图 1.1 研究内容路线图	6
图 2.1 现代主流大模型推理的两个阶段	9
图 2.2 基于因果掩码和 KV 缓存的自注意力机制	10
图 2.3 32 位浮点数向量对称量化为 8 位定点数	11
图 2.4 矩阵向量乘基本优化方法	14
图 2.5 3D 堆叠内存结构	16
图 2.6 UPMEM 硬件架构	17
图 3.1 基于查找表的矩阵向量乘	23
图 3.2 分块载入查找表卸载 GEMV 中的乘法	25
图 3.3 数制映射表卸载 GEMV 中的加法	27
图 3.4 矩阵行重排计算示意图	29
图 3.5 矩阵列重排计算示意图	32
图 4.1 FLA 指令集设计	37
图 4.2 不同 GEMV 算法查表累加语句汇编分析	38
图 4.3 基于 SIMD 指令的矩阵构建和 GEMV 计算	39
图 4.4 向量指令集设计	40
图 4.5 PIMulator 架构图	43
图 5.1 不同平台及算法下 GEMV 总吞吐	48
图 5.2 不同平台及算法下 GEMV 能效比和瓶颈分析	50
图 5.3 UPMEM 平台不同 GEMV 算法的多线程扩展性测试	52
图 5.4 UPMEM 平台不同 GEMV 算法矩阵尺寸扩展性测试	53

图 5.5 硬件架构修改对于 GEMV 算子吞吐提升——编译优化 O3	54
图 5.6 硬件架构修改对于 GEMV 算子吞吐提升——编译优化 O0	55

表目录

表 2.1 UPMEM 基本参数 19

表 3.1 FP8 常用两种格式二进制细节 26

表 4.1 软件优化性能分析 35

表 5.1 硬件平台配置 46

