

QiMeng-TensorOp: 一句提示词足以生成基于硬件原语的高性能张量算子

张续志, 彭少辉, 周其睿, 文渊博, 郭崎, 陈睿智,
朱鑫国, 熊伟强, 陈海鑫, 马聪颖, 高科, 赵琛,
武延军, 陈云霁, 李玲*

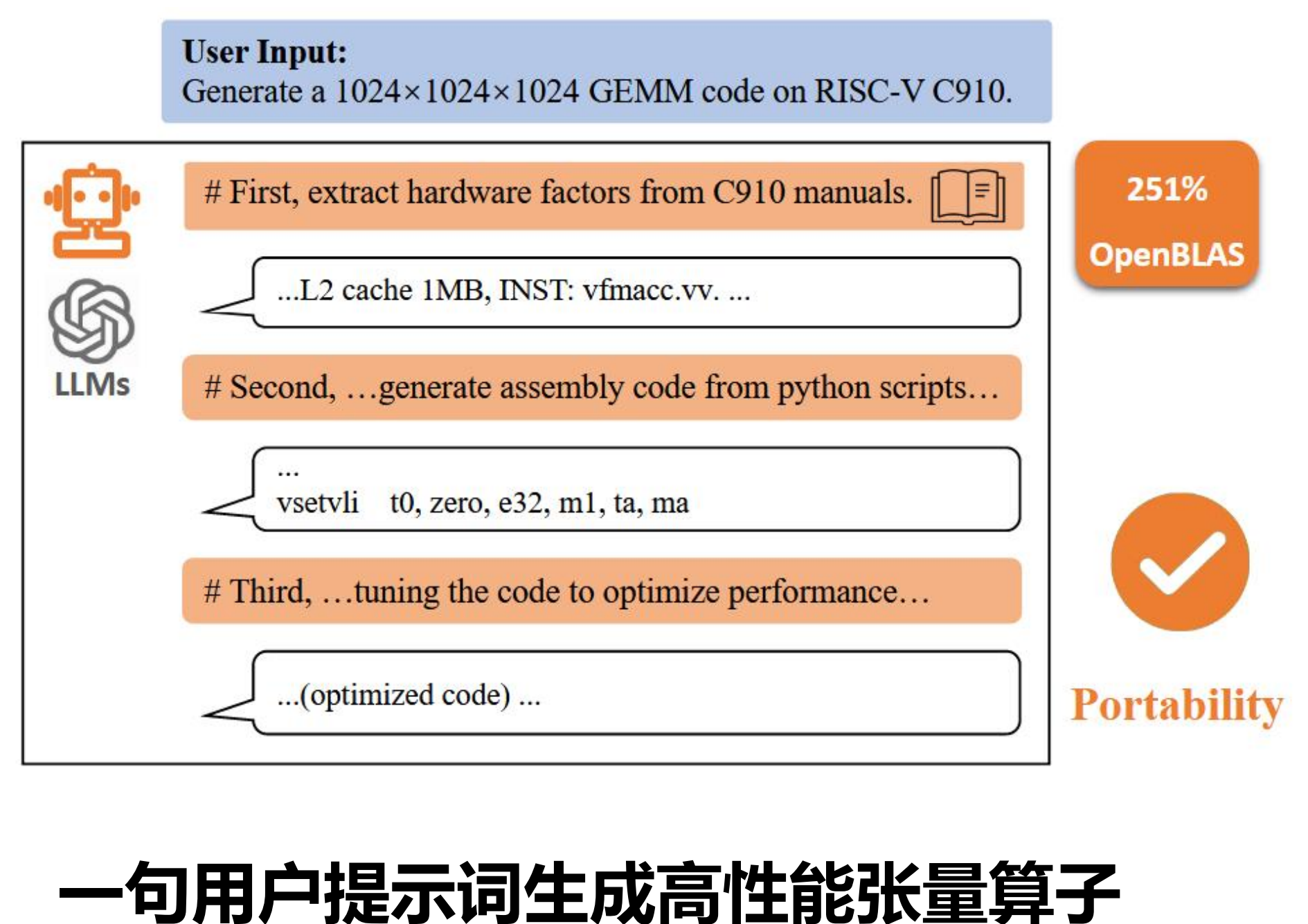
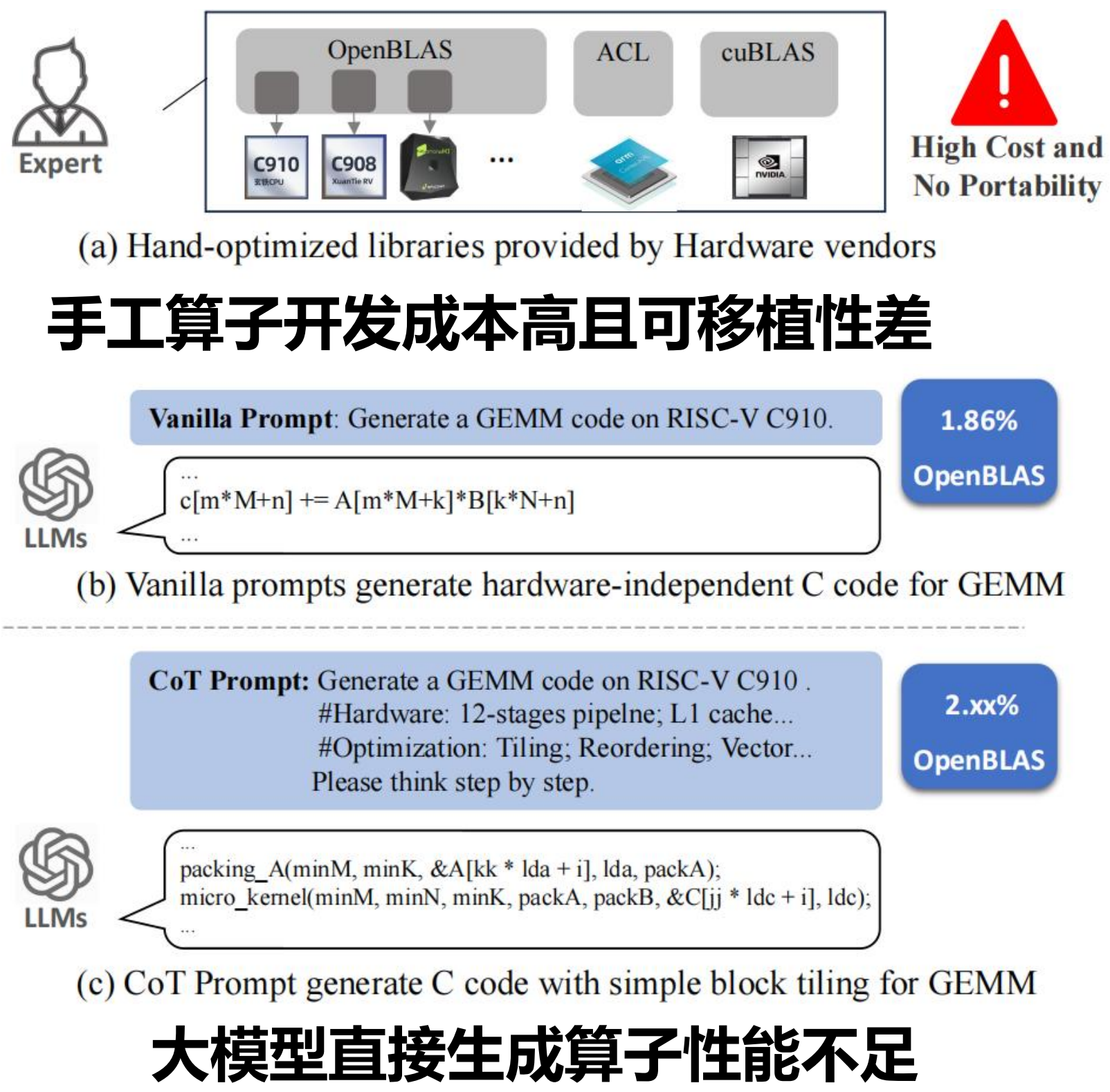
IJCAI 2025 (CCF A)

联系方式: 李玲, liling@iscas.ac.cn

工作介绍

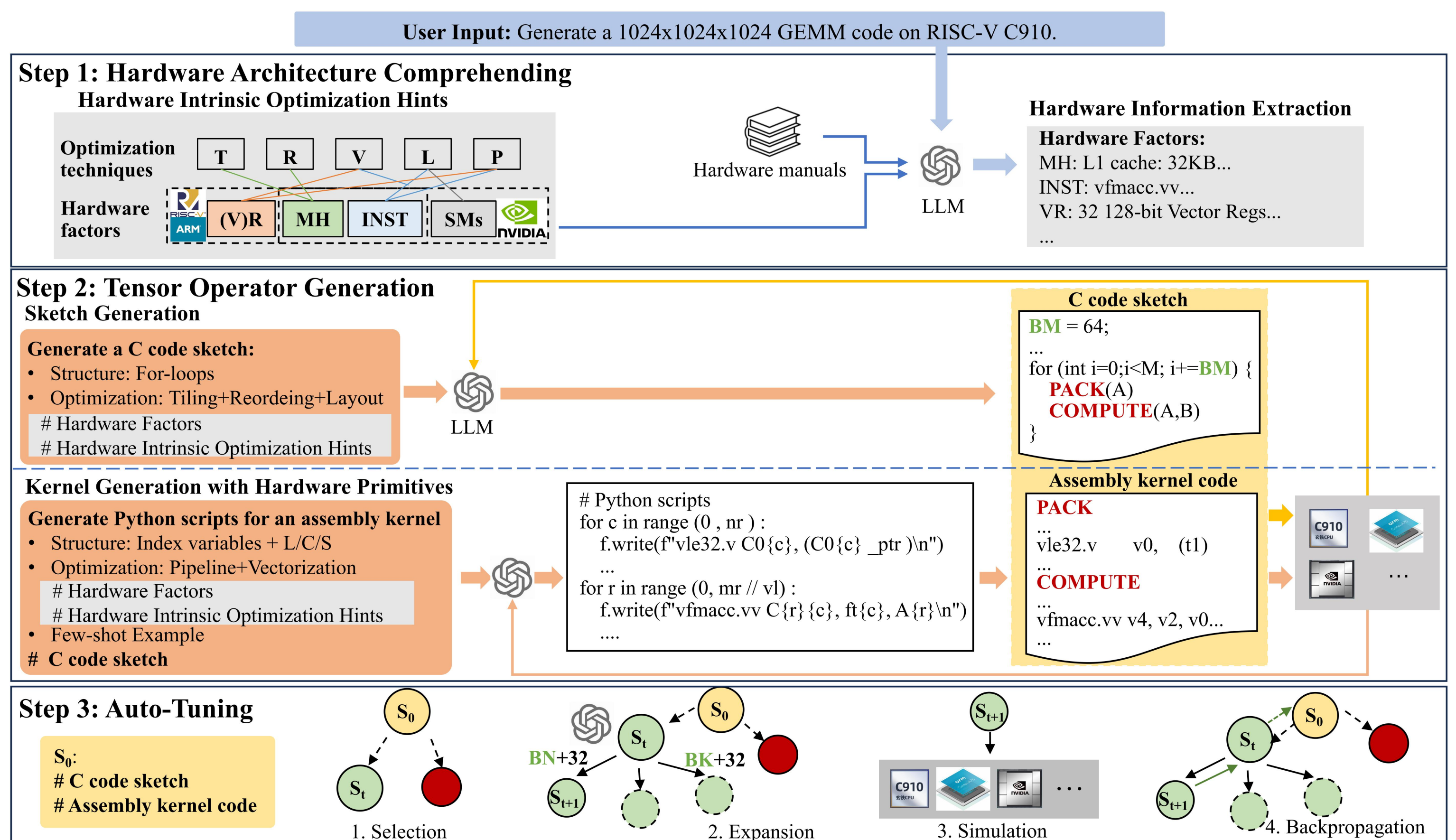
计算密集型张量算子构成了大语言模型和深度神经网络中90%的计算。由于手动优化实现不仅耗时还缺乏可移植性, 高效自动地利用硬件原语生成高性能张量算子, 对于 RISC-V、ARM 和 GPU 等多样化且不断演进的硬件架构至关重要。大语言模型擅长生成高级语言代码, 但在充分理解硬件特性并生成高性能张量算子方面仍存在不足。

我们提出了一种使用一句用户提示词自动生成张量算子的框架 (QiMeng-TensorOp), 该框架使大语言模型自动利用硬件特性, 通过硬件原语生成张量算子, 并针对不同硬件平台优化参数, 最终实现最佳性能。在多种硬件平台、前沿大语言模型及典型张量算子上的实验结果表明, QiMeng-TensorOp能有效利用各类硬件平台的计算潜力, 自动生成具有卓越性能的张量算子。



总体框架

- 硬件架构理解:** 使用优化元提示和硬件信息激发大模型对硬件架构和张量算子优化的理解
- 张量算子生成:** 通过分层策略分别生成张量算子的骨架和内核代码, 实现良好的代码结构
- 自动调优:** 利用基于大模型的MCTS发掘细微的优化机会, 进一步调优性能



实验评估

高性能且低开发代价

- 在RISC-V平台上最高达到了OpenBLAS库的**2.51倍**性能
- 在GPU平台上相比于TVM实现了显著的加速, 矩阵乘算子上最高达到**1.38倍**的性能
- 相比于高级程序员, QiMeng-TensorOp将开发代价从几天缩短到了**二十分钟**

