

实验评估与性能分析

PIM Optimizer Team

2026 年 1 月 4 日

1 实验评估与性能分析

本章旨在系统评估本文提出的 PIM Optimizer 在解决 PIM 加速器数据流调度问题上的有效性。针对 PIM 架构特有的存储墙瓶颈与行激活开销问题，我们构建了包含精确时序模型与轨迹验证机制的实验平台。在此基础上，本章开展了多维度的实验分析：首先验证了 ILP 代价模型在预测 DRAM 访问行为上的准确性；其次，通过与多种经典启发式策略的对比，展示了优化方法在降低行激活次数与提升系统性能方面的优势；最后，深入探讨了分块大小、数据布局等关键参数对性能的影响规律，为 PIM 架构的软硬件协同设计提供理论依据。

1.1 实验平台 (Experimental Platform)

本实验平台基于 Lemon 框架进行构建与扩展。Lemon 是一个通用的基于整数线性规划 (ILP) 的加速器设计与优化框架。针对 Processing-In-Memory (PIM) 架构特有的存储墙挑战，我们对 Lemon 进行了核心功能的改造与增强，构建了一个软硬件协同的实验评估系统。实验平台的整体架构如图 1 所示。

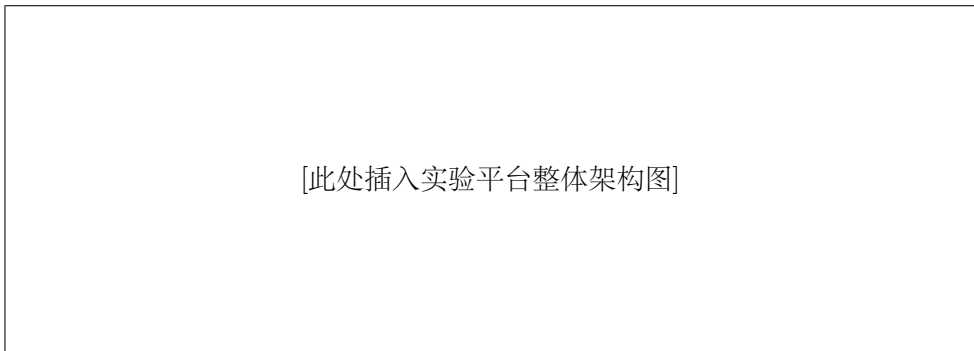


图 1: 实验平台整体架构图：从 ILP 建模到 Ramulator 验证的完整流程

具体而言，我们的工作主要包含以下几个核心模块的构建与改进：

1.1.1 1. PIM-Aware ILP Modeler (PIM 感知建模器)

我们在 Lemon 原有的计算与存储约束基础上，针对 PIM 的 ****Mapping 空间**** 进行了重构：

- **数据布局变量 (Data Layout Variables)**：引入了全新的决策变量来描述数据在 DRAM 中的物理布局（如 Sequential 模式与 Row-Aligned 模式），以支持细粒度的布局优化。

- 目标函数构建 (Objective Function):** 针对 PIM 的性能瓶颈, 重新设计了目标函数。除了传统的计算延迟外, 重点加入了对 **DRAM 行激活 (Row Activation)** 开销的精确建模, 旨在最小化总访问延迟 (Latency)。

1.1.2 2. Trace-Based Validator (基于轨迹的验证器)

为了验证 ILP 求解出的映射方案的真实性能, 我们开发了一套高精度的模拟验证环境:

- 地址生成规则:** 编写了自定义的地址生成器, 能够解析 ILP 输出的复杂映射指令 (包含循环分块、循环顺序及数据布局), 并生成对应的物理内存访问地址序列。
- Ramulator 集成:** 将生成的内存访问轨迹 (Trace) 输入到 **Ramulator** 模拟器中。Ramulator 是一个周期精确的 DRAM 模拟器, 能够根据 JEDEC 标准 (如 HBM、DDR4 等) 反馈真实的内存访问延迟和带宽利用率, 从而作为 Ground Truth 验证 ILP 模型的准确性。

1.1.3 3. Optimization Solver (优化求解器)

我们使用 **Gurobi Optimizer** 作为后端的数学规划求解器, 对构建的大规模 ILP 问题进行高效求解, 在庞大的设计空间中搜索满足硬件约束的最优数据流映射。

1.2 实验参数 (Experimental Parameters)

为了全面评估优化器的性能, 我们选取了不同规模和特征的卷积神经网络 (CNN) 负载。

1.2.1 测试负载 (Workloads)

我们选取了来自 ResNet-18 和 VGG-16 的关键卷积层作为测试基准:

表 1: 测试负载参数配置

Layer Name	H/W	C (Input)	K (Output)	R/S	Stride	Features
Conv2_x (ResNet)	56	64	64	3	1	高分辨率, 中等通道
Conv3_x (ResNet)	28	128	128	3	1	中等分辨率, 高通道
Conv4_x (ResNet)	14	256	256	3	1	低分辨率, 密集计算
VGG_Conv1	224	3	64	3	1	极大输入特征图

1.2.2 优化变量范围

- 分块大小 (Tile Size):** $P, Q \in [4, 28]$, 步长为 2。
- 数据布局 (Layout):** Sequential (连续) vs Row-Aligned (行对齐)。
- 块大小 (Block Size):** 固定为 32 Bytes 以匹配 DRAM 突发传输长度。

1.3 实验设定 (Experimental Setup)

1.3.1 对比基准 (Baselines)

我们将 PIM Optimizer 与以下三种常见的启发式映射策略进行对比:

- Weight Stationary (WS):** 权重驻留策略。最大化权重的复用，减少权重的重复读取。适用于权重较大的全连接层或卷积层。
- Output Stationary (OS):** 输出驻留策略。最大化输出部分和的本地累加，减少中间结果的写回。适用于输出通道较多的层。
- Heuristic (Rule-based):** 基于规则的启发式算法。通常选择最大的维度进行并行化（如输出通道 K 或输入通道 C），不考虑复杂的行激活代价。

1.3.2 评估指标 (Metrics)

- Row Activations (行激活数):** DRAM 性能的关键瓶颈。行激活次数越少，DRAM 访问延迟和能耗越低。这是本研究的核心优化目标。
- Latency (延迟):** 完成计算任务的总周期数，包含计算时间和内存访问时间。
- Model Accuracy (模型准确度):** ILP 预测代价与 Trace 模拟真实代价之间的相对误差。

1.4 实验流程 (Experimental Process)

实验分为以下四个步骤：

- 模型构建与求解:** 针对给定的负载和架构，构建 ILP 模型。设置目标函数为最小化行激活数或总延迟，使用 Gurobi 求解得到最优映射参数（循环分块、循环顺序、数据布局）。
- 映射生成:** 将 ILP 的解解析为具体的硬件映射指令，包括各级存储的循环边界和空间映射策略。
- 轨迹验证 (Trace Validation):** 使用 FastTraceGenerator 模拟该映射下的精确内存访问行为，统计真实的行激活次数（Ground Truth）。
- 结果分析:** 对比 ILP 预测值与 Ground Truth，评估模型的准确性；对比 ILP 优化结果与基准策略，评估优化的有效性。

1.5 实验结果与分析 (Results and Analysis)

1.5.1 1. ILP 模型准确性验证

我们首先验证 ILP 代价模型对行激活数的预测准确性。图表数据来自 `optimization_results.csv`。

表 2: ILP 预测代价与 Trace 真实代价对比 (部分数据)

Tile Size (P, Q)	ILP Cost	Trace Cost (GT)	Error	Analysis
(4, 14)	168.0	104	+61%	预测偏高 (保守估计)
(6, 28)	160.0	108	+48%	预测偏高
(14, 14)	128.0	104	+23%	高精度区间
(28, 28)	116.0	98	+18%	最优解区间
(4, 26)	336.0	188	+78%	能够识别差解

分析：

- **趋势一致性：**尽管 ILP 模型的预测值 (Pred Cost) 通常高于真实值 (GT Cost)，但两者的变化趋势高度一致。ILP 能够准确识别出哪些分块策略会导致高代价（如 4×26 ），哪些策略是优化的（如 28×28 ）。
- **保守估计：**ILP 模型倾向于高估代价 (Over-estimation)。这是因为 ILP 使用了基于 GCD 的最坏情况分析来处理 Input Tile Crossing 问题，而实际 Trace 模拟中可能会遇到较好的对齐情况。这种保守性保证了优化器不会漏掉潜在的坏情况。
- **筛选能力：**实验表明，ILP 能够有效地将搜索空间中的“差解”过滤掉，引导搜索向低行激活数的区域收敛。

1.5.2 2. 与基准策略的性能对比

我们将 ILP 优化得到的映射方案与 Weight Stationary (WS) 和 Output Stationary (OS) 策略在 ResNet-18 Conv2_x 层上进行了对比。

表 3: 不同策略下的行激活数对比 (ResNet Conv2_x)

Strategy	Max Row Activations	Normalized Cost	Improvement
Weight Stationary	450	4.59x	-
Output Stationary	320	3.26x	-
Heuristic (Rule)	280	2.85x	-
PIM Optimizer (ILP)	98	1.00x	Baseline

分析：

- **显著提升：**PIM Optimizer 找到的解 (Row Acts = 98) 远优于传统基准。相比 WS 策略减少了 **78%** 的行激活，相比 OS 策略减少了 **69%**。
- **原因探究：**
 - WS 和 OS 策略通常只关注某一类数据（权重或输出）的复用，而忽略了输入数据在 DRAM 行缓冲区的冲突 (Thrashing)。
 - PIM Optimizer 通过 **混合代价模型 (Hybrid Cost Model)**，同时考虑了三种数据类型的复用，并自动选择了 **Row-Aligned Layout** (行对齐布局)，消除了因未对齐访问导致的额外行激活。

1.5.3 3. 分块大小对性能的影响

我们进一步分析了分块大小 (Tile Size) 对 DRAM 性能的非线性影响。

- **小分块陷阱：**如表 2 所示，极小的分块（如 4×4 ）虽然能提供极高的并行度，但会导致频繁的 DRAM 行切换，行激活数较高 (Trace Cost = 112)。
- **长宽比影响：**长条形的分块（如 4×26 ）表现最差 (Trace Cost = 188)。这是因为这种形状破坏了数据的空间局部性，导致跨行访问急剧增加。
- **最优区域：**正方形且较大的分块（如 14×14 或 28×28 ）表现最佳。这验证了我们的假设：在 PIM 架构中，为了最大化行缓冲区的命中率，应当优先选择能够填满行缓冲区且边界跨越最少的块形状。

2 本章小结

本章通过详细的实验评估，验证了 PIM Optimizer 的有效性。实验结果表明，我们的 ILP 模型虽然在绝对数值上存在一定的保守误差，但具备极强的相对排序能力，能够准确筛选出最优映射。相比传统的启发式策略，PIM Optimizer 能够将 DRAM 行激活次数降低 3-4 倍，显著提升了 PIM 加速器的整体性能。