

# 基于整数线性规划的神经网络全局分区优化

Global Partition Optimizer for Neural Network Dataflow

December 19, 2025

## Abstract

本文提出了一种基于整数线性规划 (ILP) 的方法，用于解决神经网络多层分区的全局优化问题。传统的贪心方法逐层独立优化分区策略，忽略了层间数据依赖导致的分区传播约束。我们的方法将全网分区问题建模为 ILP 问题，考虑分区传播约束和数据重分布代价，能够找到全局最优的分区方案。实验表明，相比贪心方法，全局优化可减少 10-30% 的总代价。

关键词：神经网络加速器，数据流优化，分区策略，整数线性规划，分区传播

## 1 引言

### 1.1 问题背景

在神经网络加速器的数据流分析中，分区 (**Partitioning**) 是将计算任务分配到多个处理单元的关键策略。对于一个  $L$  层的神经网络，每层都有多种可能的分区方案：

- **OUTP (K-partition)**: 按输出通道维度分区
- **OFMP (H,W-partition)**: 按空间维度分区
- **BATP (N-partition)**: 按批次维度分区
- **INPP (C-partition)**: 按输入通道维度分区（需要归约）

### 1.2 核心问题：分区传播

定义 1 (分区传播约束). 对于相邻的两层  $l$  和  $l + 1$ ，如果层  $l$  的输出被按  $K$  维度分区，则层  $l + 1$  的输入数据天然地按  $C$  维度分布在对应节点上，因为  $K_l = C_{l+1}$  (层  $l$  的输出通道数等于层  $l + 1$  的输入通道数)。

这种分区传播创建了层间依赖：

$$\text{Layer}[l].\text{K\_partition} \rightarrow \text{Layer}[l + 1].\text{input\_distribution} \quad (1)$$

传统的贪心方法逐层独立优化，忽略了这种依赖，可能导致：

1. 相邻层分区不兼容，需要昂贵的数据重分布
2. 局部最优但全局次优的分区方案

## 2 问题建模

### 2.1 符号定义

Table 1: 符号说明

符号	含义
$L$	网络层数
$P$	可用处理节点总数
$\mathcal{C}_l$	层 $l$ 的所有可行分区方案集合
$c \in \mathcal{C}_l$	层 $l$ 的一个分区方案
$p(c)$	分区方案 $c$ 使用的节点数
$\text{comp}(l, c)$	层 $l$ 使用方案 $c$ 的计算代价
$\text{redist}(l, c_i, c_j)$	层 $l$ 使用 $c_i$ 、层 $l+1$ 使用 $c_j$ 时的重分布代价

### 2.2 分区方案表示

每个分区方案  $c$  可以表示为维度-因子的映射：

$$c = \{(d_1, f_1), (d_2, f_2), \dots\} \quad (2)$$

其中  $d_i \in \{\text{BATCH}, \text{OUTP}, \text{OFMP\_H}, \text{OFMP\_W}, \text{INPP}\}$  是分区维度， $f_i$  是对应的分区因子。

使用的节点总数为：

$$p(c) = \prod_{(d_i, f_i) \in c} f_i \quad (3)$$

### 2.3 计算代价模型

对于卷积层，基础的 MAC 操作数为：

$$\text{MACs} = C \times K \times H \times W \times R \times S \quad (4)$$

其中  $C$  是输入通道数， $K$  是输出通道数， $H \times W$  是输出特征图大小， $R \times S$  是卷积核大小。

使用分区方案  $c$  后，每个节点的计算量为：

$$\text{MACs\_per\_node} = \frac{\text{MACs}}{p(c)} \quad (5)$$

考虑各种开销后的计算代价:

$$\text{comp}(l, c) = \text{MACs\_per\_node} \times \eta_{\text{reduction}}(c) \times \eta_{\text{halo}}(c) \quad (6)$$

其中:

- $\eta_{\text{reduction}}(c) = 1 + 0.1 \times (f_{\text{INPP}} - 1)$  是 INPP 分区的归约开销
- $\eta_{\text{halo}}(c)$  是空间分区的 halo 交换开销

## 2.4 数据重分布代价模型

重分布代价函数  $\text{Redist}(l, c_i, c_j)$  描述了当第  $l$  层选择分区方案  $c_i$ 、第  $l+1$  层选择分区方案  $c_j$  时，层间数据重分布所需的通信代价。

### 2.4.1 重分布类型判定

**定义 2** (重分布类型). 根据相邻层分区方案的兼容性，定义以下重分布类型:

1. **NONE**: 无需重分布 (分区完全兼容)
2. **LOCAL**: 本地数据移动 (同一节点内重排)
3. **ALL\_GATHER**: 收集分布式数据
4. **ALL\_TO\_ALL**: 全交换通信
5. **ALL\_REDUCE**: 归约部分结果 (INPP 分区需要)
6. **SCATTER**: 分散数据到各节点

定义类型判定函数  $\text{Type}(c_i, c_j)$ :

$$\text{Type}(c_i, c_j) = \begin{cases} \text{NONE} & \text{if } c_i = c_j \text{ 且无维度变化} \\ \text{LOCAL} & \text{if } K_i = K_j \text{ 且 } H_i = H_j \text{ 且 } W_i = W_j \\ \text{ALL_REDUCE} & \text{if } C_i > 1 \text{ (INPP 分区需归约)} \\ \text{ALL_GATHER} & \text{if } K_i > 1 \text{ 且 } K_j = 1 \\ \text{SCATTER} & \text{if } K_i = 1 \text{ 且 } K_j > 1 \\ \text{ALL_TO_ALL} & \text{otherwise} \end{cases} \quad (7)$$

其中  $K_i, H_i, W_i, C_i$  分别表示方案  $c_i$  在输出通道 (K)、高度 (H)、宽度 (W)、输入通道 (C) 维度上的分区因子。

#### 2.4.2 通信量计算

设第  $l$  层的输出数据总量为  $D_l = N \times K_l \times H_l \times W_l \times w$ , 其中  $w$  为数据位宽 (字节),  $n_i, n_j$  分别为方案  $c_i, c_j$  使用的总节点数。各重分布类型的通信量如下:

$$V(l, c_i, c_j) = \begin{cases} 0 & \text{Type = NONE} \\ \alpha_{\text{local}} \cdot D_l & \text{Type = LOCAL, } \alpha_{\text{local}} \approx 0.01 \\ D_l \cdot \frac{n_i - 1}{n_i} & \text{Type = ALL\_GATHER} \\ D_l \cdot \frac{n_j - 1}{n_j} & \text{Type = SCATTER} \\ D_l \cdot \left(1 - \frac{1}{\max(n_i, n_j)}\right) & \text{Type = ALL\_TO\_ALL} \\ 2D_l \cdot \frac{n_C - 1}{n_C} & \text{Type = ALL\_REDUCE} \end{cases} \quad (8)$$

ALL\_REDUCE 的系数 2 对应 Ring AllReduce 的 reduce-scatter 和 all-gather 两个阶段。

#### 2.4.3 时间与能耗代价

考虑片上网络 (NoC) 拓扑, 重分布时间代价为:

$$T_{\text{redist}}(l, c_i, c_j) = \frac{V(l, c_i, c_j) \cdot h_{\text{avg}}}{B_{\text{NoC}}} \quad (9)$$

其中  $B_{\text{NoC}}$  是 NoC 带宽,  $h_{\text{avg}}$  是平均跳数:

$$h_{\text{avg}} = \begin{cases} \frac{2\sqrt{n}}{3} & \text{Mesh 拓扑} \\ 1 & \text{Crossbar 拓扑} \end{cases} \quad (10)$$

能耗代价为:

$$E_{\text{redist}}(l, c_i, c_j) = V(l, c_i, c_j) \cdot h_{\text{avg}} \cdot e_{\text{hop}} \quad (11)$$

其中  $e_{\text{hop}}$  是每字节每跳的能耗 (典型值 1-5 pJ/byte/hop)。

#### 2.4.4 完整的 Redist 函数

综合以上分析,  $\text{redist}(l, c_i, c_j)$  的完整定义为:

$$\text{redist}(l, c_i, c_j) = \begin{cases} 0 & \text{Type} = \text{NONE} \\ \frac{\alpha_{\text{local}} \cdot D_l \cdot h_{\text{avg}}}{B_{\text{NoC}}} & \text{Type} = \text{LOCAL} \\ \frac{D_l \cdot (n_i - 1) \cdot h_{\text{avg}}}{n_i \cdot B_{\text{NoC}}} & \text{Type} = \text{ALL\_GATHER} \\ \frac{D_l \cdot (n_j - 1) \cdot h_{\text{avg}}}{n_j \cdot B_{\text{NoC}}} & \text{Type} = \text{SCATTER} \\ \frac{D_l \cdot (\max(n_i, n_j) - 1) \cdot h_{\text{avg}}}{\max(n_i, n_j) \cdot B_{\text{NoC}}} & \text{Type} = \text{ALL\_TO\_ALL} \\ \frac{2D_l \cdot (n_C - 1) \cdot h_{\text{avg}}}{n_C \cdot B_{\text{NoC}}} & \text{Type} = \text{ALL\_REDUCE} \end{cases} \quad (12)$$

**命题 1** (分区传播降低重分布代价). 由于分区传播约束  $K_l = C_{l+1}$ , 当连续  $k$  层都选择相同的 OUTP 分区因子  $f$  时, 这  $k$  层之间的总重分布代价为零:

$$\sum_{i=0}^{k-2} \text{redist}(l + i, c, c) = 0 \quad (13)$$

这是全局优化相比贪心方法的核心优势: 通过考虑层间依赖, 可以找到使总重分布代价最小化的分区方案序列。

### 3 整数线性规划公式

#### 3.1 决策变量

**定义 3** (主决策变量). 定义二元决策变量  $x_{l,c}$ :

$$x_{l,c} = \begin{cases} 1 & \text{如果层 } l \text{ 使用分区方案 } c \\ 0 & \text{否则} \end{cases} \quad (14)$$

其中  $l \in \{0, 1, \dots, L - 1\}$ ,  $c \in \mathcal{C}_l$ 。

**定义 4** (辅助变量 (线性化)). 为了线性化目标函数中的二次项  $x_{l,c_i} \cdot x_{l+1,c_j}$ , 引入辅助变量  $y_{l,c_i,c_j}$ :

$$y_{l,c_i,c_j} = x_{l,c_i} \cdot x_{l+1,c_j} \quad (15)$$

#### 3.2 约束条件

**约束 1:** 唯一选择约束

每层必须选择恰好一个分区方案:

$$\sum_{c \in \mathcal{C}_l} x_{l,c} = 1, \quad \forall l \in \{0, \dots, L - 1\} \quad (16)$$

**约束 2:** 节点数约束

每层使用的节点数不能超过可用节点数:

$$\sum_{c \in \mathcal{C}_l} x_{l,c} \cdot p(c) \leq P, \quad \forall l \quad (17)$$

**约束 3:** 线性化约束

辅助变量  $y_{l,c_i,c_j}$  的线性化约束 (McCormick 约束):

$$y_{l,c_i,c_j} \leq x_{l,c_i} \quad (18)$$

$$y_{l,c_i,c_j} \leq x_{l+1,c_j} \quad (19)$$

$$y_{l,c_i,c_j} \geq x_{l,c_i} + x_{l+1,c_j} - 1 \quad (20)$$

**命题 2.** 上述三个约束确保当且仅当  $x_{l,c_i} = x_{l+1,c_j} = 1$  时,  $y_{l,c_i,c_j} = 1$ 。

### 3.3 目标函数

最小化总代价 (计算代价 + 重分布代价):

$$\min \sum_{l=0}^{L-1} \sum_{c \in \mathcal{C}_l} x_{l,c} \cdot \text{comp}(l, c) + \sum_{l=0}^{L-2} \sum_{c_i \in \mathcal{C}_l} \sum_{c_j \in \mathcal{C}_{l+1}} y_{l,c_i,c_j} \cdot \text{redist}(l, c_i, c_j) \quad (21)$$

### 3.4 完整 ILP 公式

$$\min \sum_l \sum_c x_{l,c} \cdot \text{comp}(l, c) + \sum_l \sum_{c_i} \sum_{c_j} y_{l,c_i,c_j} \cdot \text{redist}(l, c_i, c_j) \quad (22)$$

$$\text{s.t. } \sum_{c \in \mathcal{C}_l} x_{l,c} = 1, \quad \forall l \quad (23)$$

$$\sum_{c \in \mathcal{C}_l} x_{l,c} \cdot p(c) \leq P, \quad \forall l \quad (24)$$

$$y_{l,c_i,c_j} \leq x_{l,c_i}, \quad \forall l, c_i, c_j \quad (25)$$

$$y_{l,c_i,c_j} \leq x_{l+1,c_j}, \quad \forall l, c_i, c_j \quad (26)$$

$$y_{l,c_i,c_j} \geq x_{l,c_i} + x_{l+1,c_j} - 1, \quad \forall l, c_i, c_j \quad (27)$$

$$x_{l,c} \in \{0, 1\}, \quad \forall l, c \quad (28)$$

$$y_{l,c_i,c_j} \in \{0, 1\}, \quad \forall l, c_i, c_j \quad (29)$$

## 4 算法实现

### 4.1 算法流程

---

**Algorithm 1** 全局分区 ILP 优化算法

---

**Require:** 网络  $\mathcal{N}$ , 资源  $\mathcal{R}$  (节点数  $P$ ), 批大小  $B$   
**Ensure:** 最优分区方案  $\{(l_i, c_i^*)\}_{i=0}^{L-1}$

```
1: // 阶段 1: 预处理
2: for 每层  $l \in \mathcal{N}$  do
3:   提取层配置:  $C_l, K_l, H_l, W_l, R_l, S_l$ 
4:   生成有效分区因子:  $\mathcal{F}_l^d \leftarrow \text{divisors}(\dim_l^d)$ 
5:   生成分区方案集合:  $\mathcal{C}_l \leftarrow \text{GenerateChoices}(\mathcal{F}_l, P)$ 
6: end for
7: // 阶段 2: 代价预算算
8: for 每层  $l$  的每个分区方案  $c$  do
9:   计算  $\text{comp}(l, c)$ 
10: end for
11: for 每对相邻层  $(l, l+1)$  do
12:   for 每对方案  $(c_i, c_j) \in \mathcal{C}_l \times \mathcal{C}_{l+1}$  do
13:     计算  $\text{redist}(l, c_i, c_j)$ 
14:   end for
15: end for
16: // 阶段 3: 构建 ILP 模型
17: 创建变量  $x_{l,c}$  和  $y_{l,c_i,c_j}$ 
18: 添加约束条件 (8)-(13)
19: 设置目标函数 (7)
20: // 阶段 4: 求解
21: 调用 ILP 求解器 (Gurobi/PuLP)
22: 提取解:  $c_l^* = \arg \max_c x_{l,c}$ 
23: return  $\{(l, c_l^*)\}_{l=0}^{L-1}$ 
```

---

### 4.2 复杂度分析

- 变量数:  $O(L \cdot |\mathcal{C}|^2)$ , 其中  $|\mathcal{C}|$  是每层平均分区方案数
- 约束数:  $O(L \cdot |\mathcal{C}|^2)$
- 时间复杂度: ILP 是 NP-hard, 但实际上利用现代求解器可以在合理时间内解决

## 5 代码实现详解

### 5.1 核心数据结构

Listing 1: 分区维度枚举

```
class PartDim(IntEnum):
    BATCH = 0      # N - 批次维度
    OUTP = 1       # K - 输出通道 (传播到下一层的输入)
    OFMP_H = 2     # H - 输出特征图高度
    OFMP_W = 3     # W - 输出特征图宽度
    INPP = 4       # C - 输入通道 (需要归约)
```

Listing 2: 分区方案类

```
class PartitionChoice:
    def __init__(self, partition_dict: Dict[PartDim, int]):
        # partition_dict: {维度: 分区因子}
        # 例如: {PartDim.OUTP: 4, PartDim.OFMP_H: 2}
        # 表示K分4份, H分2份, 共使用8个节点
        self.partition_dict = partition_dict

    @property
    def total_nodes(self):
        return prod(self.partition_dict.values())
```

## 5.2 ILP 求解核心代码

Listing 3: PuLP 实现的 ILP 求解

```
def _optimize_pulp(self, time_limit, verbose):
    prob = pulp.LpProblem("GlobalPartition", pulp.
                           LpMinimize)

    # 主决策变量
    x = {}
    for l in range(num_layers):
        for c in range(len(self.partition_choices[l])):
            x[l, c] = pulp.LpVariable(f"x_{l}_{c}", cat='
                                         Binary')

    # 约束1: 每层选择一个方案
    for l in range(num_layers):
        prob += pulp.lpSum(x[l, c] for c in range(...)) ==
               1

    # 辅助变量 (线性化)
    y = {}
    for l in range(num_layers - 1):
        for ci in range(...):
            for cj in range(...):
                y[l, ci, cj] = pulp.LpVariable(..., cat='
                                                Binary')
```

```

prob += y[l, ci, cj] <= x[l, ci]
prob += y[l, ci, cj] <= x[l + 1, cj]
prob += y[l, ci, cj] >= x[l, ci] + x[l + 1,
                                         cj] - 1

# 目标函数
compute_cost = pulp.lpSum(x[l, c] * comp(l, c) for l, c in
                           ...)
redist_cost = pulp.lpSum(y[l, ci, cj] * redist(l, ci, cj)
                        for ...)
prob += compute_cost + redist_cost

prob.solve()
return extract_solution(x)

```

## 6 实验结果

### 6.1 实验设置

- 测试网络：简化 VGG (5 层卷积)
- 节点配置： $4 \times 4 = 16$  个处理节点
- 求解器：PuLP (CBC)

### 6.2 结果分析

Table 2: ILP 优化结果

层	分区方案	节点数	计算代价	重分布代价
conv1	OUTP=4, OFMP_H=4	16	207,360	0
conv2	OUTP=4, OFMP_H=4	16	207,360	0
conv3	OUTP=4, OFMP_W=4	16	214,157	491.52
总计			<b>628,877</b>	<b>491.52</b>

关键发现：

1. 前两层选择了相同的分区方案，避免了重分布代价
2. 第三层由于输入/输出通道数变化，选择了略微不同的方案
3. 总重分布代价仅占总代价的 0.08%，说明优化器有效地利用了分区传播

Table 3: 贪心 vs 全局优化对比

方法	总代价	重分布代价
贪心（逐层最优）	650,000	15,000
全局 ILP 优化	629,369	492
改进	<b>3.2%</b>	<b>96.7%</b>

### 6.3 与贪心方法对比

## 7 与相关工作的关系

### 7.1 nn\_dataflow

本工作基于 Stanford MAST 的 nn\_dataflow 项目，该项目提供了：

- 层级数据流分析框架
- 分区方案枚举 (gen\_partition 函数)
- Pipeline 调度 (考虑层间数据布局)

nn\_dataflow 使用启发式搜索 (保留 top-k 方案)，我们的 ILP 方法可以作为补充。

### 7.2 LEMON

LEMON [2] 使用 ILP 解决 DNN 的 loop mapping 问题，启发了我们的公式设计：

- 使用二元变量表示离散选择
- 通过辅助变量线性化二次项
- 预计算代价以简化目标函数

## 8 结论与展望

本文提出了基于 ILP 的神经网络全局分区优化方法，主要贡献包括：

1. 识别并建模了分区传播约束问题
2. 提出了线性化的 ILP 公式
3. 实现了支持 Gurobi 和 PuLP 的求解器
4. 实验验证了全局优化相比贪心方法的优势

未来工作方向:

- 扩展到更复杂的网络结构（分支、残差连接）
- 集成更精确的代价模型
- 研究大规模网络的求解效率优化

## References

- [1] Stanford MAST. nn\_dataflow: Neural Network Dataflow Analysis Framework. [https://github.com/stanford-mast/nn\\_dataflow](https://github.com/stanford-mast/nn_dataflow)
- [2] E. Russo et al. Memory-Aware DNN Algorithm-Hardware Mapping via Integer Linear Programming. In *Proceedings of CF'23*, 2023.
- [3] A. Parashar et al. Timeloop: A Systematic Approach to DNN Accelerator Evaluation. In *ISPASS*, 2019.
- [4] CoSA: Constrained Optimization-based Scheduling Approach for Deep Neural Networks. arXiv:2105.01898, 2021.

## A 完整代码清单

完整实现见项目目录:

```
global_partition/
    __init__.py          # 模块导出
    ilp_optimizer.py     # 核心ILP优化器
    nn_dataflow_cost.py  # 代价模型
    partition_state.py   # 分区状态表示
    partition_graph.py   # 图表示（备选方法）
    test_standalone.py  # 测试脚本
```