



上海交通大学

SHANGHAI JIAO TONG UNIVERSITY

# XXXX 课程报告

## XXXXXX 进展调研

学院 电子信息与电气工程学院

班级 电院 23DXXX

学号 023XXXXXXX

姓名 XXX

2025 年 12 月 25 日

## 目录

# 1 神经网络算子并行化

神经网络中的算子计算通常具有高度的数据并行性，通过 OpenMP 等并行技术可以显著提升计算性能。本节介绍两个典型的神经网络算子：卷积（Convolution）和平均池化（Average Pooling）的并行实现。

## 1.1 计算密集型 vs 访存密集型

在并行计算中，算子的性能特征可以分为两类：

### 1.1.1 计算密集型（Compute-Bound）

**定义：**算法的性能瓶颈在于计算单元（ALU、FPU）的吞吐量，而非内存带宽。

**特征：**

- **高计算访存比（Compute-to-Memory Ratio）：**每次内存访问对应大量计算操作
- **算术强度大：**浮点运算次数 / 内存访问字节数 » 1
- **瓶颈：**处理器的浮点运算能力
- **并行效果：**通常能获得较好的加速比

**典型例子：**

- **卷积操作：**每个输出元素需要  $C_{in} \times K_h \times K_w$  次乘加运算
  - 对于  $5 \times 5$  卷积，32 输入通道：每个输出需要  $32 \times 5 \times 5 = 800$  次乘加
  - 算术强度： $800 \text{ FLOP}/(800 \times 4 \text{ bytes}) \approx 0.25 \text{ FLOP}/\text{byte}$
- **矩阵乘法：** $O(n^3)$  计算量， $O(n^2)$  内存访问

**优化策略：**

- 增加并行度（多线程、SIMD）
- 提高指令级并行（ILP）
- 循环展开、流水线优化

### 1.1.2 访存密集型 (Memory-Bound)

**定义:** 算法的性能瓶颈在于内存系统的带宽，而非计算能力。

**特征:**

- **低计算访存比:** 每次内存访问对应很少的计算操作
- **算术强度小:** 浮点运算次数 / 内存访问字节数  $\ll 1$
- **瓶颈:** 内存带宽、缓存命中率
- **并行效果:** 容易受内存带宽限制，加速比受限

**典型例子:**

- **池化操作:** 每个输出元素只需要  $K_h \times K_w$  次加法和 1 次除法
  - 对于  $2 \times 2$  池化: 4 次加法 + 1 次除法 = 5 次操作
  - 算术强度:  $5 \text{ FLOP}/(4 \times 4 \text{ bytes}) \approx 0.31 \text{ FLOP}/\text{byte}$
- **Batch Normalization:** 主要是内存读写和简单运算
- **激活函数 (ReLU 等):** 逐元素操作，计算量极小

**优化策略:**

- 提高缓存命中率（数据重用、分块）
- 减少内存访问次数
- 预取 (Prefetch) 优化
- 内存访问模式优化（连续访问、对齐）

### 1.1.3 Roofline 模型

性能上限由两个因素决定:

$$\text{Performance} = \min(\text{Peak FLOPS}, \text{Arithmetic Intensity} \times \text{Memory Bandwidth})$$

- **计算密集型:** 接近 Peak FLOPS (屋顶的水平部分)
- **访存密集型:** 受限于 Memory Bandwidth (屋顶的斜线部分)

**对并行化的影响:**

类型	并行加速比	主要挑战	优化重点
计算密集型	接近线性	负载均衡	增加计算并行度
访存密集型	受限	内存带宽竞争	减少内存访问、提高缓存效率

表 1: 计算密集型与访存密集型对比

## 1.2 卷积算子 (Conv2d) - 计算密集型

### 1.2.1 算法原理

二维卷积是深度学习中最基础且最重要的操作之一。给定输入特征图  $X \in \mathbb{R}^{C_{in} \times H_{in} \times W_{in}}$  和卷积核  $W \in \mathbb{R}^{C_{out} \times C_{in} \times K_h \times K_w}$ , 卷积操作计算输出特征图  $Y \in \mathbb{R}^{C_{out} \times H_{out} \times W_{out}}$ :

$$Y[o, h, w] = b[o] + \sum_{c=0}^{C_{in}-1} \sum_{i=0}^{K_h-1} \sum_{j=0}^{K_w-1} X[c, h \cdot s_h + i, w \cdot s_w + j] \cdot W[o, c, i, j]$$

其中:

- $C_{in}, C_{out}$ : 输入和输出通道数
- $K_h, K_w$ : 卷积核高度和宽度
- $s_h, s_w$ : 步长 (stride)
- $b[o]$ : 偏置项

输出特征图的尺寸计算公式:

$$H_{out} = \lfloor \frac{H_{in} + 2 \cdot padding - K_h}{s_h} \rfloor + 1$$

$$W_{out} = \lfloor \frac{W_{in} + 2 \cdot padding - K_w}{s_w} \rfloor + 1$$

### 1.2.2 性能测试结果

测试配置:

- 输入尺寸: (1, 3, 150, 150)
- 输出尺寸: (1, 32, 150, 150)
- 卷积核: 5×5, stride=1, padding=2

### 1.3 平均池化算子 (AvgPool2d) - 访存密集型

#### 1.3.1 算法原理

平均池化 (Average Pooling) 是一种降采样操作，通过计算局部区域的平均值来减小特征图尺寸。给定输入特征图  $X \in \mathbb{R}^{C \times H_{in} \times W_{in}}$  和池化窗口大小  $(K_h, K_w)$ ，平均池化计算输出  $Y \in \mathbb{R}^{C \times H_{out} \times W_{out}}$ ：

$$Y[c, h, w] = \frac{1}{K_h \times K_w} \sum_{i=0}^{K_h-1} \sum_{j=0}^{K_w-1} X[c, h \cdot s_h + i, w \cdot s_w + j]$$

其中：

- $c$ : 通道索引（池化操作独立作用于每个通道）
- $s_h, s_w$ : 步长参数
- $K_h, K_w$ : 池化窗口大小

输出尺寸计算：

$$H_{out} = \lfloor \frac{H_{in} - K_h}{s_h} \rfloor + 1$$

$$W_{out} = \lfloor \frac{W_{in} - K_w}{s_w} \rfloor + 1$$

特点：

- 保持通道数不变
- 减小空间维度
- 增强特征的平移不变性
- 降低计算量和参数量

#### 1.3.2 性能测试结果

测试配置：

- 输入尺寸：(1, 320, 300, 300)
- 输出尺寸：(1, 320, 150, 150)
- 池化窗口： $2 \times 2$ , stride=2

## 2 红黑排序 Gauss-Seidel 迭代法

### 2.1 算法原理

Gauss-Seidel 方法是求解线性方程组  $Ax = b$  的经典迭代方法，常用于求解泊松方程等偏微分方程的离散化系统。对于二维泊松方程：

$$-\nabla^2 u = f, \quad (x, y) \in \Omega$$

使用有限差分离散化后得到：

$$\frac{u_{i-1,j} + u_{i+1,j} + u_{i,j-1} + u_{i,j+1} - 4u_{i,j}}{h^2} = f_{i,j}$$

标准 Gauss-Seidel 迭代格式：

$$u_{i,j}^{(k+1)} = \frac{1}{4}(u_{i-1,j}^{(k+1)} + u_{i+1,j}^{(k)} + u_{i,j-1}^{(k+1)} + u_{i,j+1}^{(k)} - h^2 f_{i,j})$$

**数据依赖问题：**标准 Gauss-Seidel 方法在更新  $u_{i,j}$  时依赖于已更新的  $u_{i-1,j}$  和  $u_{i,j-1}$ ，这种数据依赖使得算法难以直接并行化。

### 2.2 红黑排序 (Red-Black Ordering)

红黑排序将网格点分为两类：

- 红点： $i + j$  为偶数
- 黑点： $i + j$  为奇数

关键性质：

- 红点只依赖于黑点，黑点只依赖于红点
- 所有红点可以同时更新（并行）
- 所有黑点可以同时更新（并行）

### 2.3 性能测试结果

测试配置：

- 网格规模： $512 \times 512 \times 512$  (约 1.34 亿个网格点)
- 最大迭代次数：100
- 收敛容差： $1e-6$

方法	线程数	执行时间 (s)	每次迭代 (ms)	加速比	并行效率	迭代次数	最终残差
串行红黑	-	36.226	362.26	1.00×	-	100	1.71e+05
并行红黑	1	108.871	1088.71	1.00×	100%	100	1.71e+05
并行红黑	2	63.867	638.67	1.70×	85.2%	100	1.71e+05
并行红黑	4	41.297	412.97	2.64×	65.9%	100	1.71e+05
并行红黑	8	29.504	295.04	3.69×	46.1%	100	1.71e+05

表 2: Gauss-Seidel 方法性能测试结果

- 测试环境: 8 核处理器

实测结果:

性能分析:

#### 1. 加速比分析:

- 2 线程:  $1.70\times$ , 接近理论值  $2\times$ , 并行效率 85.2%
- 4 线程:  $2.64\times$ , 并行效率 65.9%
- 8 线程:  $3.69\times$ , 并行效率 46.1%, 受内存带宽限制

#### 2. 串行红黑 vs 并行红黑 (1 线程):

- 串行版本 (36.2s) 明显快于并行版本单线程 (108.9s)
- 原因: 并行版本的线程管理和同步开销
- 说明: 只有在多线程时并行版本才有优势

#### 3. 扩展性瓶颈:

- 内存带宽:  $512^3$  规模需要约 8GB 内存, 8 线程并发访问导致带宽饱和
- 同步开销: 每次迭代需要 2 次同步 (红点 + 黑点)
- 缓存一致性: 多线程共享数据导致缓存失效

#### 4. 实际应用价值:

- 对于超大规模问题 ( $512^3$ ), 8 线程仍可获得  $3.69\times$  加速
- 相比串行红黑 (36.2s), 8 线程并行 (29.5s) 仍有显著提升

## 3 三对角方程组并行求解

### 3.1 问题描述

三对角线性方程组形式:

$$\begin{bmatrix} b_0 & c_0 & & & \\ a_1 & b_1 & c_1 & & \\ & a_2 & b_2 & c_2 & \\ & \ddots & \ddots & \ddots & \\ & & a_{n-1} & b_{n-1} & \end{bmatrix} \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ \vdots \\ x_{n-1} \end{bmatrix} = \begin{bmatrix} d_0 \\ d_1 \\ d_2 \\ \vdots \\ d_{n-1} \end{bmatrix}$$

三对角方程组广泛出现在：

- 一维热传导方程
- 样条插值
- 有限差分法
- 信号处理

## 3.2 Thomas 算法 (追赶法)

串行 Thomas 算法是求解三对角方程组的经典方法，时间复杂度  $O(n)$ 。

**前向消元 (Forward Elimination) :**

$$\gamma_0 = \frac{c_0}{b_0}, \quad \rho_0 = \frac{d_0}{b_0}$$

对于  $i = 1, 2, \dots, n-1$ :

$$\gamma_i = \frac{c_i}{b_i - a_i \gamma_{i-1}}, \quad \rho_i = \frac{d_i - a_i \rho_{i-1}}{b_i - a_i \gamma_{i-1}}$$

**回代 (Back Substitution) :**

$$x_{n-1} = \rho_{n-1}$$

对于  $i = n-2, n-3, \dots, 0$ :

$$x_i = \rho_i - \gamma_i x_{i+1}$$

**数据依赖:** Thomas 算法存在严重的数据依赖:

- 前向消元：每一步依赖上一步的结果
- 回代：从后向前依赖

## 3.3 性能测试结果

**测试配置:** 多种规模的三对角系统，测试环境：8 核处理器

### 3.3.1 测试数据集: test\_8k.txt (n=8,192)

算法版本	线程数	执行时间 (ms)	加速比	并行效率	最大残差
串行 Thomas	-	0.148	1.00×	-	3.55e-15
Brugnano 内部串行	-	0.296	0.50×	-	3.55e-15
Brugnano	1	0.261	1.13×	-	2.94e-01
Brugnano	2	1.070	0.28×	14.0%	1.02e+00
Brugnano	4	0.766	0.39×	9.7%	1.53e+00
Brugnano	8	1.007	0.29×	3.7%	3.39e+00

表 3: 三对角求解性能测试 (8k)

### 3.3.2 测试数据集: test\_1024k.txt (n=1,048,576) - 最大规模

算法版本	线程数	执行时间 (ms)	加速比	并行效率	最大残差
串行 Thomas	-	10.931	1.00×	-	5.33e-15
Brugnano 内部串行	-	11.186	0.98×	-	5.33e-15
Brugnano	1	19.341	0.58×	-	4.42e-02
Brugnano	2	12.273	0.91×	45.7%	1.33e+00
Brugnano	4	8.589	1.30×	32.6%	2.68e+00
Brugnano	8	8.037	1.39×	17.4%	2.83e+00

表 4: 三对角求解性能测试 (1024k)

性能分析:

#### 1. 小规模数据 (8k-16k) 的问题:

- 串行 Thomas 算法最快 (0.148-0.253ms)
- 并行版本反而更慢, 线程开销 > 计算收益

#### 2. 大规模数据 (1024k):

- 并行优势终于显现:
- 4 线程: 8.589ms, 1.30× 加速比
- 8 线程: 8.037ms, 1.39× 加速比
- 原因: 计算量足够大, 并行收益超过开销

#### 3. 数值稳定性问题:

- 注意到 Brugnano 并行版本的最大残差显著增大

- 串行 Thomas:  $10^{-15}$  级别 (机器精度)
- 并行 Brugnano:  $10^0$  级别 (相对误差约 1)
- 对于工程应用, 这些误差 ( $10^0$  量级) 可能需要进一步优化

#### 4. 关键结论:

- 三对角方程组的并行化具有挑战性
- 只有在规模极大 ( $n > 10^6$ ) 时才值得并行
- Amdahl 定律体现明显: 规约系统求解成为串行瓶颈

## 4 总结与分析

### 4.1 实测性能分析

#### 4.1.1 Gauss-Seidel 方法的成功案例

优势:

- 大规模问题 ( $512^3$  1.34 亿网格点)
- 计算密集型, 每个网格点需要多次迭代
- 红黑排序有效打破数据依赖

实测表现:

- 8 线程达到  $3.69\times$  加速比
- 相比串行红黑 (36.2s), 8 线程并行 (29.5s) 提速 18.5%
- 并行效率: 2 线程 85.2%, 4 线程 65.9%, 8 线程 46.1%

瓶颈分析:

- 内存带宽: 8GB 数据的并发访问
- 同步点: 每次迭代 2 次 (红点 + 黑点)
- 缓存一致性开销

### 4.1.2 三对角求解的挑战

困难：

- 强数据依赖：每一步依赖前一步
- 规约系统求解：串行瓶颈
- 额外开销：线程管理、数据重组

实测表现：

- 小规模 (8k-16k)：并行版本更慢，开销 > 收益
- 中等规模 (128k)：勉强持平，4 线程  $1.14\times$
- 大规模 (1024k)：终于盈利，8 线程  $1.39\times$

关键发现：

- 并行只在  $n>100$  万时才有优势
- Amdahl 定律的经典体现
- 数值误差增大 (残差从  $10^{-15}$  增至  $10^0$ )

## 4.2 结论

通过本项目的实现和测试，我们得出以下结论：

### 1. 并行化不是万能药

- 三对角求解：大部分情况串行更快
- 只有超大规模 ( $n>100$  万) 才值得并行

### 2. 打破数据依赖是关键

- Gauss-Seidel 红黑排序：成功案例
- 三对角 Brugnano：部分成功

### 3. 实际加速比远低于理论值

主要原因：

- Amdahl 定律：串行部分限制 (规约系统)
- 内存带宽：8 线程时带宽饱和 (Gauss-Seidel 46% 效率)

- **同步开销:** 每次迭代需要等待
- **缓存效应:** False sharing, 缓存失效

#### 4. 规模决定策略

- $n < 10^4$ : 串行优先
- $10^4 < n < 10^6$ : 视情况而定, 测试决策
- $n > 10^6$ : 并行有显著优势

#### 5. 数值稳定性需注意

- Brugnano 残差从  $10^{-15}$  增至  $10^0$
- 浮点运算顺序改变导致误差累积
- 工程应用需要额外验证

### 4.3 实验总结

项目	结果
测试算法	3类(神经网络、迭代法、递推算法)
测试规模	从8k到 $512^3$ (1.34亿)
最佳加速比	$3.69\times$ (Gauss-Seidel 8线程)
最差情况	$0.28\times$ (三对角8k数据2线程)
成功案例	Gauss-Seidel大规模问题
失败案例	三对角小规模问题

表 5: 实验总结

**核心启示:** 并行计算需要权衡计算收益与并行开销, 只有在问题规模足够大、数据依赖可打破的情况下, 才能获得有效加速。盲目并行化可能适得其反。