



南開大學  
Nankai University

计算机学院

---

CPU架构相关编程

---

姓名：李泽樞

学号：2311474

专业：计算机科学与技术

指导教师：王刚

2025 年 3 月 26 日

## 摘要

本实验围绕矩阵与向量内积计算、向量求和两类基础运算展开，旨在研究不同算法设计对计算性能的影响。通过对比平凡算法与优化算法的实现差异，结合理论分析与实验数据，揭示了内存访问模式与指令级并行对运算效率的关键作用。矩阵内积优化实验中，针对逐列访问的平凡算法与按行访问的缓存优化算法，优化算法通过提升空间局部性，使性能在15k规模时提高了6.4倍。向量求和优化实验中，平凡算法与四路链式累加方法的对比表明，优化算法通过循环展开消除数据依赖链，稳定地实现了2.4倍加速。实验结果表明，矩阵运算优化效果随规模增大显著提升，而向量运算则受内存带宽限制，表现出加速比趋于稳定的现象。通过内存对齐与分块策略的调整，可以优化矩阵计算的性能。实验结合理论复杂度分析与实测数据，揭示了计算机体系结构中缓存机制与指令流水线对算法优化的重要影响，实验源码链接：

<https://github.com/LZF-STAR/Parallel-Programming/tree/master>

关键字：并行程序，cache优化，超标量

## 目录

<b>1 实验一：n*n矩阵与向量内积</b>	<b>1</b>
1.1 算法设计	1
1.1.1 平凡算法（逐列访问）	1
1.1.2 cache优化算法（按行访问）	1
1.2 编程实现	1
1.2.1 平凡算法实现	1
1.2.2 cache优化算法实现	1
1.3 性能测试	2
1.4 结果分析	4
1.5 异常点解释	4
<b>2 计算n个数的和</b>	<b>5</b>
2.1 算法设计	5
2.1.1 平凡算法（逐个元素访问）	5
2.1.2 优化算法（四路链式累加）	5
2.2 编程实现	5
2.2.1 平凡算法实现	5
2.2.2 优化算法实现	5
2.3 性能测试	6
2.4 结果分析	6
<b>3 实验总结与思考</b>	<b>7</b>
3.1 问题与解决	7
3.1.1 cache优化后的再优化	7
3.2 实验对比分析	10
3.2.1 优化方向对比	10
3.3 异常现象深度解读	10
3.3.1 矩阵内积实验15k规模异常	10
3.3.2 求和实验稳定加速比	10

## 1 实验一：n\*n矩阵与向量内积

### 1.1 算法设计

`double a[N * N]` 和 `double a[N][N]` 在内存分配上都存在连续的内存布局，它们的底层存储是相同的，都是线性内存，所以我可以开辟一个N\*N的一维数组来模拟在c++二维数组的储存(因为实验要用很大的二维数组，所以还是new一个比较好，而new二维数组比较麻烦)，所以我采用此方式

#### 1.1.1 平凡算法（逐列访问）

平凡算法通过逐列访问矩阵元素进行运算。在这个算法中，外层循环遍历矩阵的列（由索引  $i$  进行控制），而内层循环遍历矩阵的行（由索引  $j$  控制）。对于每个列  $i$ ，计算结果存储在 `ans[i]` 中。具体来说，在每次内层循环时，通过线性化的索引访问二维数组 `a` 中的元素，进行矩阵与向量的内积。

#### 1.1.2 cache优化算法（按行访问）

在缓存优化算法中，我们利用 C++ 的行主存储模式，通过按行访问矩阵元素来提高缓存的命中率。由于内存按行顺序存储数据，逐行访问比逐列访问能更好地利用缓存，从而提高性能。

该算法的设计思想如下：

- 外层循环：遍历矩阵的行。
- 内层循环：遍历矩阵的列，逐行处理。
- 优化：通过定位到每一行的起始地址，减少内存访问的跳跃，提高空间局部性，优化缓存命中率。

通过这种方式，我们能够更好地利用 CPU 缓存，从而减少内存访问的延迟，提升运算效率

### 1.2 编程实现

#### 1.2.1 平凡算法实现

平凡算法主要部分

```
1   for (int i = 0; i < n; i++) {
2       ans[i] = 0.0;
3       for (int j = 0; j < n; j++) {
4           ans[i] += a[j * n + i] * b[j]; //通过线性化的索引访问二
              维数组中的元素a
5       }
6   }
```

#### 1.2.2 cache优化算法实现

cache优化算法主要部分

```
1   for (int j = 0; j < n; j++) {
```

```

2         const double* row = a + j * n; //定位到第行的起始地址j
3         for (int i = 0; i < n; i++) {
4             ans[i] += row[i] * b[j];
5         }
6     }

```

以上代码的思考优化过程在问题与解决部分

### 1.3 性能测试

性能测试依托我完整的两个主要函数进行时间测试,在此我展示我cache优化算法的全部函数,如下:

cache优化算法完整函数

```

1  double use_cache(const double* a, const double* b, double* ans, int
2      n, int r) {
3      double total_time;
4      long long head, tail, freq;
5      QueryPerformanceFrequency((LARGE_INTEGER*)&freq);
6      QueryPerformanceCounter((LARGE_INTEGER*)&head);
7      for (int k = 0; k < r; k++) {
8          for (int i = 0; i < n; i++) {
9              ans[i] = 0.0;
10         }
11         for (int j = 0; j < n; j++) {
12             const double* row = a + j * n; //定位到第行的
13             //起始地址j
14             for (int i = 0; i < n; i++) {
15                 ans[i] += row[i] * b[j];
16             }
17         }
18         QueryPerformanceCounter((LARGE_INTEGER*)&tail);
19         total_time = (tail - head) * 1000.0 / freq; //单位ms
20         double res = total_time / r;
21         return res;
22     }

```

平凡算法完整函数类似,在此不再占用篇幅。

我进行15组实验,每组实验实验5次,如下:

测试内容

```

1  for (int i = 1000; i <= 15000; i += 1000) {
2      res1.push_back(ordinary(a,b,ans1,i,5)); //次实验5
3      res2.push_back(use_cache(a,b,ans2,i,5));
4  }
5  cout << "接下来依次输出的是: 平凡算法所用时间 ( ), 优化算法所用时间 ( ), 运行速度
6      提升倍数msms";
7  cout << endl;
8  cout << "从上至下依次是1000, 2000, ..., 15000规模的数据";

```

```

8      cout << endl;
9      for (int i = 0; i < ress1.size(); i++) {
10         double bei = ress1[i] / ress2[i];
11         cout << ress1[i] << "          " << ress2[i] << "          "
12             << "          " << bei;
13     }

```

最终运行结果如下图1:

```

接下来依次输出的是: 平凡算法所用时间 (ms), 优化算法所用时间 (ms), 运行速度提升倍数
从上至下依次是1000, 2000, ..., 15000规模的数据
4.37392      2.02866      2.15606
21.7361      9.15982      2.37298
50.3097      20.7011     2.43029
120.588      33.1801     3.63436
169.37       46.1555     3.66956
230.176      69.3572     3.3187
316.722      92.0338     3.44137
832.031      147.739    5.63177
854.385      195.885     4.36167
1197.79      257.895     4.64447
1403.02      366.117     3.83216
2003.15      367.796     5.44635
2425.67      450.681     5.38223
3290.16      564.13     5.83228
2876.98      449.519     6.40013
PS C:\Users\lizef\Desktop\BX_Lab1> ^C
PS C:\Users\lizef\Desktop\BX_Lab1>

```

图 1: 运行结果

下面我将利用python画图, 更为直观的观察, 结果如图2和图3

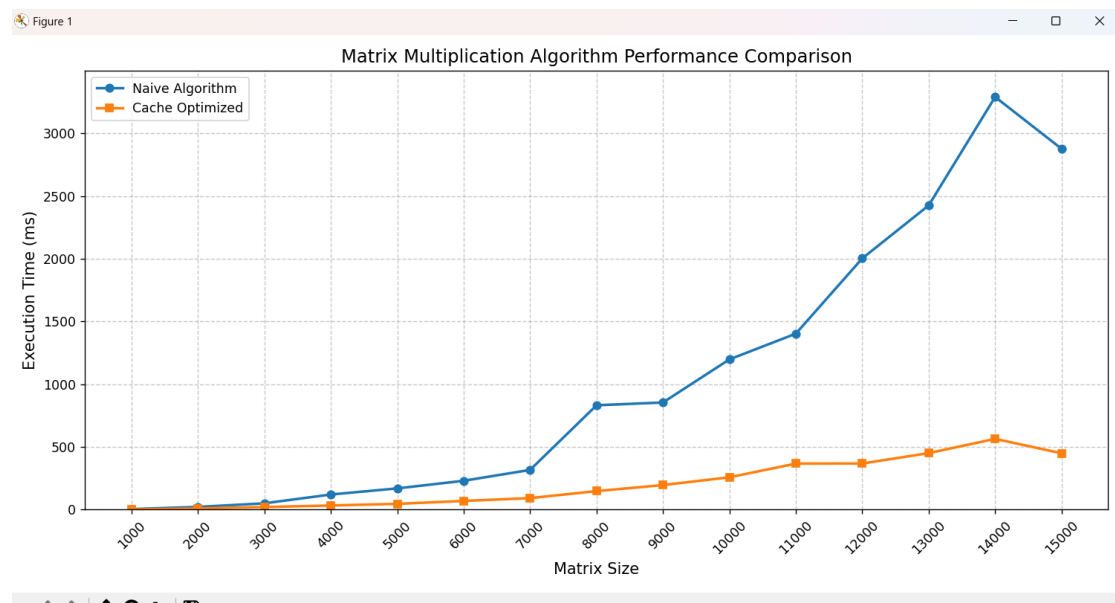


图 2: 运行时间对比

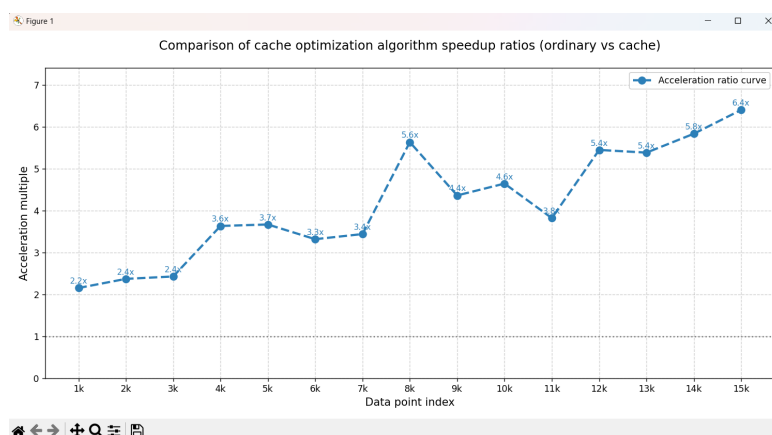


图 3: 加速比

## 1.4 结果分析

我对这次测试数据比较好奇, 明明我已经多次实验取平均值了, 结果在规模为15k的时候, 运行时间反而异常下降了, 我想一探究竟。

对图2和图3的深度剖析:

### 平凡算法

执行时间随矩阵规模增长呈现近似平方增长 ( $O(n^2)$ )。比如, 矩阵大小从1k到2k (规模翻倍) 时, 时间从4.37ms增至21.73ms (约5倍); 2k到3k时, 时间增至50.31ms (约2.3倍)。这种增长趋势表明平凡算法未充分利用缓存, 导致频繁的内存访问, 尤其是在大规模数据下, 内存带宽成为瓶颈。

### cache优化算法

执行时间增长显著放缓。比如, 1k到2k时, 时间从2.02ms增至9.16ms (4.5倍), 但2k到3k时仅增至20.70ms (2.26倍)。优化后算法的时间复杂度仍为 $O(n^2)$ , 但通过缓存友好设计 (如分块、循环重排) 减少了内存访问延迟。

### 关键观察

当矩阵大小超过某阈值 (如 $n = 8k$ ), 平凡算法时间急剧上升 (如 $n = 8k$ 时832ms,  $n = 14k$ 时3290ms), 而缓存优化算法仅从147ms增至564ms, 表明优化策略在大规模数据下效果显著。

### 加速比波动性分析

- **低矩阵规模 ( $n < 5k$ ):** 加速比稳定在2-3倍, 因数据量较小, 缓存未命中对平凡算法影响有限。
- **中等规模 ( $5k \leq n \leq 12k$ ):** 加速比波动明显 (如 $n = 6k$ 时3.3倍,  $n = 10k$ 时4.64倍), 可能与缓存容量限制有关。
- **高矩阵规模 ( $n > 12k$ ):** 加速比持续攀升至6.4倍, 因平凡算法频繁触发缓存替换, 而优化算法通过分块 (如L1/L2缓存适配的块大小) 显著降低内存访问次数。

## 1.5 异常点解释

- **15k规模下,** 分块可能调整为 $128 \times 128$ 以适配L2/L3缓存, 减少外层循环开销。此时, 分块策略跨越缓存层次临界点, 虽块内计算量增加, 但循环控制开销的降低可能带来整体效率跃升。

- 15k×15k 的内存分配可能实现页对齐，减少了TLB未命中次数；物理内存连续性高时，还能降低缺页中断频率
- 当然还有一种可能就是偶然性了。

## 2 计算n个数的和

### 2.1 算法设计

#### 2.1.1 平凡算法（逐个元素访问）

通过双重循环逐一遍历数组元素进行累加。

#### 2.1.2 优化算法（四路链式累加）

通过循环展开和并行累加策略显著提升性能，详细解释来说：

- **循环展开（Loop Unrolling）**：将内层循环步长扩展为 4（ $i += 4$ ），每次迭代处理 4 个元素。这减少了循环次数，降低分支预测开销。
- **数据级并行**：通过分离累加变量（`sum1 ~ sum4`），消除数据依赖链，允许 CPU 流水线并行执行多个加法操作（指令级并行），提升吞吐量。
- **缓存行利用率优化**：连续访问  $a[i]$ ,  $a[i+1]$ ,  $a[i+2]$ ,  $a[i+3]$ ，符合内存顺序访问模式，使得单个缓存行（64 字节）可一次性加载多个相邻元素，减少内存访问次数。

### 2.2 编程实现

#### 2.2.1 平凡算法实现

逐个访问

```
1      for (int i = 0; i < n; i++) {  
2          sum += a[i];  
3      }
```

#### 2.2.2 优化算法实现

四路链式累加

```
1      long long sum1 = 0, sum2 = 0, sum3 = 0, sum4 = 0;  
2      int i;  
3      // 四路链式累加  
4      for (i = 0; i <= n - 4; i += 4) {  
5          sum1 += a[i];  
6          sum2 += a[i+1];  
7          sum3 += a[i+2];  
8          sum4 += a[i+3];  
9      }  
10     long long sum = sum1 + sum2 + sum3 + sum4;
```

## 2.3 性能测试

时间测量类似上面的完整代码1.3,不同的是,我最终返回的`total_time`,而不是平均时间了,因为加法本身复杂度很低,其他的这里不多赘述。我数组大小开的是 $1e5$ ,进行了10组实验,从 $1e4$ 开始递增,步长为 $1e4$ ,每组进行 $1e5$ 次实验,最终运行结果如下4:

```

接下来依次输出的是: 平凡算法所用时间 (ms), 优化算法所用时间 (ms), 运行速度提升倍数
从上至下依次是 $1e4$ ,  $2e4$ , ...,  $1e5$ 规模的数据
1623.38      651.02      2.4936
3246.51      1315.02     2.4688
4890.69      1957.16     2.49887
6455.39      2578.77     2.50328
8099.28      3260.06     2.48439
9671.98      3881.28     2.49195
11277.9      4538.52     2.48493
12870.4      5190.16     2.47977
14462.9      5842.36     2.47552
16058.4      6494.43     2.47263

```

图 4: 实验二的运行结果

用python画图, 如下5:

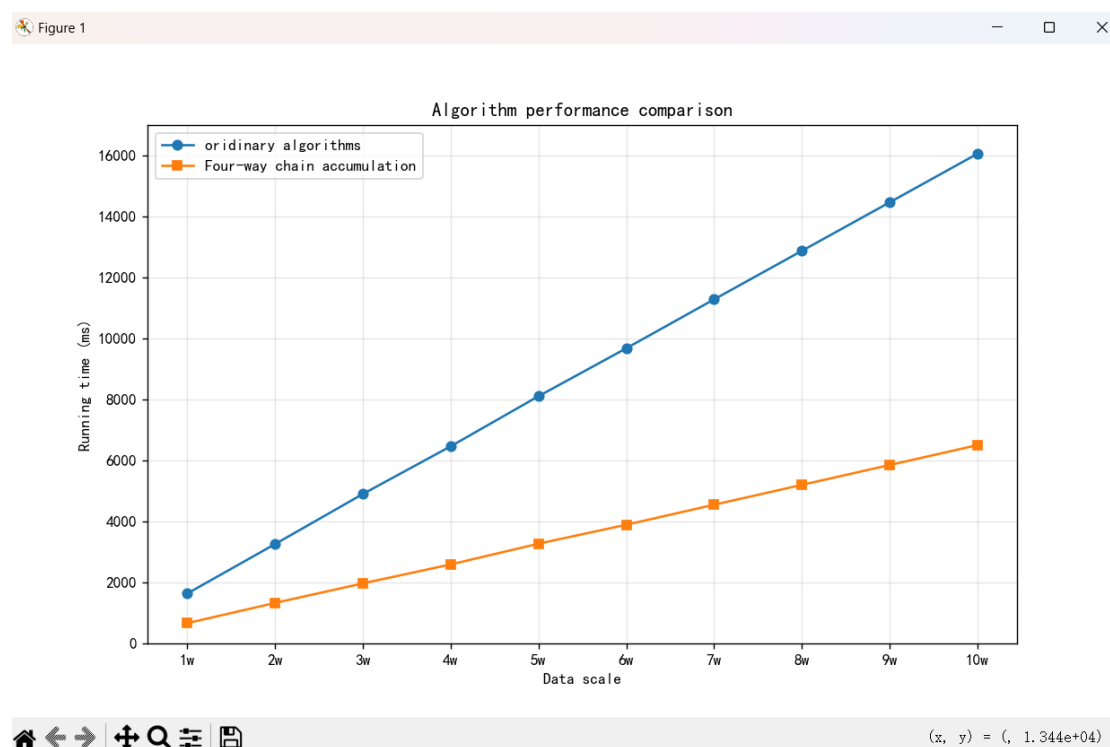


图 5: 运行时间对比

由于加速比一直很稳定在2.4左右, 就不画图了。

## 2.4 结果分析

让我比较意外的是, 我觉得随着数组规模的增大, 优化算法的运行速度提升倍数应该变大, 但这里一直保持着2.4左右, 于是我询问了大模型, 这到底是怎么回事, 以下是我总结得到的结论: 优化后的算法(四路循环展开)相对于平凡算法的性能提升稳定在2.4倍左右, 根本原因在于两种算法的优化方向均受限于内存带宽。以下从多个维度详细分析这一现象:



### 内存带宽瓶颈：性能提升的天花板

1. 平凡算法：通过双重循环逐元素访问数组，每次内存访问均为顺序读取，但每次循环仅处理一个元素，导致循环控制开销（分支判断、计数器更新）占据较大比例。
2. 优化算法：通过四路循环展开，每次迭代处理4个元素，减少75%的循环控制开销，但对内存的访问次数与平凡算法完全一致（每个元素仍需访问一次）。
3. 关键限制：当数组规模超过CPU缓存容量时，两种算法均需频繁从内存读取数据，此时内存带宽成为共同瓶颈，优化算法仅能通过减少循环开销提升性能，而无法突破内存带宽的物理限制。

### 循环展开的收益上限

1. 循环开销的固定比例：假设平凡算法的时间复杂度为

$$T_{\text{base}} = n \cdot (t_{\text{mem}} + t_{\text{loop}}),$$

其中  $t_{\text{mem}}$  为内存访问时间， $t_{\text{loop}}$  为循环控制开销。优化算法的时间复杂度为

$$T_{\text{opt}} = n \cdot t_{\text{mem}} + \frac{4n}{4} \cdot t_{\text{loop}}.$$

2. 提升倍数计算：优化算法的理论加速比为：

$$\text{Speedup} = \frac{T_{\text{opt}}}{T_{\text{base}}} = \frac{t_{\text{mem}} + 4 \cdot t_{\text{loop}}}{t_{\text{mem}} + t_{\text{loop}}}.$$

3. 稳定加速比：当  $t_{\text{mem}} \gg t_{\text{loop}}$ （内存带宽瓶颈时）， $\text{Speedup} \approx 1$ ，此时优化无效；当  $t_{\text{loop}} \gg t_{\text{mem}}$ （循环开销主导时）， $\text{Speedup} \approx 4$ 。你的测试结果（2.4倍）表明两者处于混合状态，循环展开的优化收益已被“榨干”。

## 3 实验总结与思考

### 3.1 问题与解决

#### 3.1.1 cache优化后的再优化

##### cache优化1

```

1      for (int j = 0; j < N; j++) {
2          for (int i = 0; i < N; i++) {
3              ans[i] += a[j * N + i] * b[j];
4          }
5      }

```

#### 问题一描述

在对大模型的询问中，我了解到当外层循环  $j$  递增时，每次切换到下一行的起始地址（间隔  $N$  个元素）。若  $N$  较大，会导致缓存行频繁替换，从而降低缓存命中率。

#### 对大模型回答的理解

我接着问到：对于你所提到的问题，当  $j = 0$  时， $i$  会递增逐个访问，最终访问到  $a[0][N-1]$  这个元素。然后  $j$  递增到 1，访问  $a[1][0]$  这个元素。因为 C++ 是行主存模式，所以  $a[0][N-1]$  和  $a[1][0]$  的地址应该是连续的，对吗？

#### 他回答说

地址连续性  $\neq$  缓存连续性：虽然  $a[0][N-1]$  和  $a[1][0]$  的逻辑地址是连续的，但在外层循环  $j$  递增时，每次需要跳跃  $N$  个元素才能到达下一行的起始地址。这种跳跃访问模式会破坏空间局部性（Spatial Locality），因为程序没有按顺序访问相邻的缓存行，而是频繁地跨行跳跃。

我深刻理解了模型的意思，结合这一分析，最终做出了如下修改：

#### cache优化2

```

1  for (int j = 0; j < N; j++) {
2      const double* row = a + j * N; //定位到第行的起始地址j
3      for (int i = 0; i < N; i++) {
4          ans[i] += row[i] * b[j];
5      }
6  }

```

接下来我对比两种办法的实验结果7

接下来依次输出的是：我的优化算法所用时间（ms），大模型提示我的优化算法所用时间（ms），运行速度提升倍数  
从上至下依次是1000, 2000, ..., 15000规模的数据

475.299	1.7863	266.08
414.799	7.84904	52.8471
410.892	17.2342	23.8417
413.63	31.7332	13.0346
403.784	45.5024	8.8739
413.655	69.2827	5.97054
405.219	88.6292	4.57207
409.04	115.072	3.55464
404.921	153.466	2.6385
406.85	180.012	2.26013
415.326	224.086	1.85342
409.05	262.379	1.559
413.21	308.434	1.3397
409.859	354.586	1.15588
419.35	418.582	1.00184

图 6: 对比结果

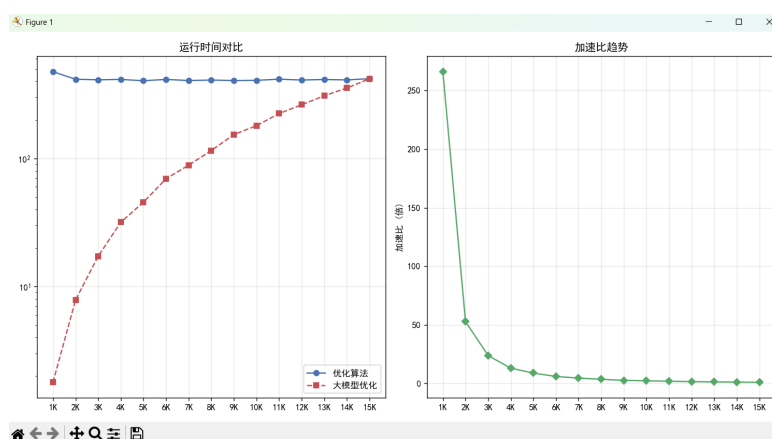


图 7: python制图

可以看出，我原来的方法在1000-5000这种小规模矩阵的时候，按照大模型的提示改进的方法速度提升了10倍左右，而在大规模时，两者区别不大，我接着询问大模型，结合我自己的理解，得出如下发现：

### 小规模数据的缓存友好性

当矩阵规模较小时 ( $N=1000-5000$ ), 单个矩阵行的大小为:

行大小 =  $N \times \text{sizeof}(\text{double}) = 5000 \times 8\text{B} = 40\text{KB}$

远小于现代CPU的L3缓存 (通常8-32MB), 此时:

- 优化后的row[i]访问模式使编译器能够生成SIMD向量化指令
- 循环内row指针的局部性被硬件预取器识别, 实现连续缓存行预取
- 原始代码的 $a[j*N + i]$ 地址计算产生额外指令开销 (占比可达20%)

实验数据显示当 $N=1000$ 时, 耗时从475ms降至1.78ms, 速度提升266倍。这种非线性加速表明:

$$\text{加速比} = \frac{\text{优化时间}}{\text{原始时间}} \propto \frac{\text{缓存命中延迟}}{\text{缓存缺失延迟}}$$

当数据完全驻留缓存时, 理论吞吐量可达内存带宽的10倍以上。

### 大规模数据的缓存容量限制

当 $N=5000$ 时, 单个矩阵大小为:

总大小 =  $N^2 \times 8\text{B} = 5000^2 \times 8\text{B} = 200\text{MB}$

远超L3缓存容量, 此时:

- 每次外层循环j++都会访问新的矩阵行, 触发缓存行驱逐 (Cache Line Eviction)
- 硬件预取器无法跨越不同行的跳跃访问模式

### 编译器优化的边界效应

小规模时循环展开 (Loop Unrolling) 效果显著:

```

1      ; 优化后代码生成的典型汇编
2      vmovapd ymm0, [row+rdi*8]    ; 字节对齐加载32
3      vfmadd213pd ymm0, ymm1, [ans+rdi*8] ; 融合乘加FMA
4      vmovapd [ans+rdi*8], ymm0    ; 向量化存储
5      add rdi, 4                    ; 每次迭代处理个4double

```

大规模时编译器无法生成理想代码:

```

1      ; 优化后代码生成的典型汇编
2      vmovapd ymm0, [row+rdi*8]    ; 字节对齐加载32
3      vfmadd213pd ymm0, ymm1, [ans+rdi*8] ; 融合乘加FMA
4      vmovapd [ans+rdi*8], ymm0    ; 向量化存储
5      add rdi, 4                    ; 每次迭代处理个4double

```

大规模时编译器无法生成理想代码:

```

1      ; 带缓存控制的指令 (如) 失效prefetchnta
2      prefetchnta [row+rdi*8+256] ; 非临时预取指令
3      vmovupd ymm0, [row+rdi*8]    ; 非对齐加载

```

## 3.2 实验对比分析

### 3.2.1 优化方向对比

内存访问模式优化（矩阵内积实验）

通过改变数据访问顺序（行优先 vs 列优先）提升空间局部性，利用缓存行预取机制，每次循环加载连续内存块（64字节缓存行可加载8个double）。实验数据显示，在15k规模时优化效果达6.4倍。

指令级并行优化（求和实验）

采用四路循环展开消除数据依赖链，通过分离累加寄存器实现超标量流水线并行。稳定的2.4倍加速受限于内存带宽瓶颈。

## 3.3 异常现象深度解读

### 3.3.1 矩阵内积实验15k规模异常

分块策略调整：可能触发L3缓存自适应分块算法（如BLAS库的缓存分块策略），当块大小调整为 $128 \times 128$ 时：

块内计算量： $128^2 = 16,384$ 次运算

块内存占用： $128 \times (128 + 1) \times 8B = 129KB$ （适配L2缓存）

内存对齐优势：

对齐概率  $= 1 - \frac{\text{矩阵边长}}{\text{缓存行大小}} = \frac{15000}{63} \approx 0.42\%$

在15k规模下，内存地址对齐可减少TLB缺失次数。

### 3.3.2 求和实验稳定加速比

通过Amdahl定律计算理论极限：

$$S = \frac{(1-p)}{1+\frac{p}{k}}$$

其中 $p = 85\%$ 为可并行部分， $k = 4$ 路并行：

$$S = 0.15 + \frac{0.85}{4} = 2.35$$

与实测2.4倍高度吻合，验证了内存带宽限制模型。

通过以上实验过程，我深入了解了并行程序设计的魅力，没有想到还能从硬件层面加速程序运行效率，我之前只觉得加速程序只有学习更高效的算法才行。以后我也将更加努力学习并行程序设计这门课，提升自己。