

PJ4 优化

Your classmate posted a new Note.

PJ4 优化

PJ4 优化

PJ4是关于 [Iterative Stencil Loops](#)中2D Jacobi iteration的计算。从题设的评判标准看，迭代的正确性检验，是通过取你的结果与倒数第一步与倒数第二步差值的最小值平方，再累加与EPS=1e-5相比较。

本题的题设是，给定大小为\$N\$的矩阵，应用2D Jacobi iteration计算矩阵中非边缘元素的\$k\$步迭代。注意到2D Jacobi iteration的计算模式，是取矩阵非边缘位置的相邻元素算平均值。

我们的第一步优化可以是：

- 奇数步：从第0行第0列开始计数，在奇数行\$m\$计算第\$m\$行第\$(1, 3, 5, 7, \dots)\$列的值，在偶数行\$n\$计算第\$n\$行第\$(2, 4, 6, 8, \dots)\$列的值。跳过第0行、第\$N-1\$行、第0列，第\$N-1\$列；
- 偶数步：从第0行第0列开始计数，在奇数行\$m\$计算第\$m\$行第\$(2, 4, 6, 8, \dots)\$列的值，在偶数行\$n\$计算第\$n\$行第\$(1, 3, 5, 7, \dots)\$列的值。跳过第0行、第\$N-1\$行、第0列，第\$N-1\$列。

正确性留待读者自证。

在结束了算法上的优化后，注意到我们每次以间隔列元素的方式计算，不方便使用SIMD优化，cache访问也不友好。因此我们将矩阵进行重排。

- 首先，考虑使用aligned_alloc(32, <size>)存放重排后的矩阵。这个函数可以在分配内存时按第一个参数规定的大小进行内存对齐，从而减少访问开销。函数的具体用法见[c++reference](#)。
- 其次，考虑到我们的计算方法，我们使用两组矩阵\$A, B\$，分别存放奇数步的计算结果与偶数步的计算结果。设原矩阵为\$P\$，则\$A, B\$与原矩阵的对应关系是：
$$A = \{a_{\lfloor i/2 \rfloor, \lfloor j/2 \rfloor} : a_{\lfloor i/2 \rfloor, \lfloor j/2 \rfloor} = P_{\lfloor i \rfloor, \lfloor i \rfloor + \lfloor j \rfloor \pmod{2}} \equiv 1, i, j \in [0, N) \cap \text{mathbb{Z}} \}$$
$$B = \{b_{\lfloor i/2 \rfloor, \lfloor j/2 \rfloor} : b_{\lfloor i/2 \rfloor, \lfloor j/2 \rfloor} = P_{\lfloor i \rfloor, \lfloor i \rfloor + \lfloor j \rfloor \pmod{2}} \equiv 0, i, j \in [0, N) \cap \text{mathbb{Z}} \}$$

在矩阵的重排过程中，注意到当\$N\$分别为奇数和偶数时，\$A\$和\$B\$每行存放的列元素数量有所区别。具体来说，以\$A\$为例。

- 当\$N\$为偶数时，\$A\$每行存放的列元素均为\$N/2\$个；
- 当\$N\$为奇数时，\$A\$的奇数行存放\$\lfloor N/2 \rfloor + 1\$个元素，\$A\$的偶数行存放\$\lfloor N/2 \rfloor\$个元素。
- \$B\$的情况与之类似，请读者自行推导。

虽然\$N\$分别为偶数和奇数时，\$A\$和\$B\$的存放模式不同，我们仍然可以以相同的模式进行计算（见下）。

- 通过以上对应关系，
- 在奇数步时，我们从\$A\$矩阵提取数据，进行计算，并将结果存放至\$B\$矩阵。计算公式如下：
$$B_{\lfloor i \rfloor, \lfloor j \rfloor} = (A_{\lfloor i-1 \rfloor, \lfloor j \rfloor} + A_{\lfloor i+1 \rfloor, \lfloor j \rfloor} + A_{\lfloor i, j-1 \rfloor} + A_{\lfloor i, j+1 \rfloor}) / 4, i \in [1, N-1], j \in [1 - (\text{and } 1), \lfloor N/2 \rfloor]$$
- 在偶数步时，我们从\$B\$矩阵提取数据，进行计算，并将结果存放至\$A\$矩阵。计算公式如下：
$$A_{\lfloor i \rfloor, \lfloor j \rfloor} = (B_{\lfloor i-1 \rfloor, \lfloor j \rfloor} + B_{\lfloor i, j \rfloor} + B_{\lfloor i+1 \rfloor, \lfloor j \rfloor} + B_{\lfloor i, j+1 \rfloor}) / 4, i \in [1, N-1], j \in [i \pmod{2}, \lfloor N/2 \rfloor]$$

由于目前矩阵的访问模式已经变成了按照列元素逐次访问，因此我们可以引入SIMD。

- \$SIMD\$计算上述表达式略。注意到AVX指令集最多支持同时计算4个double，因此我们选择_mm256类指令。需要注意的是，在[Intrinsics Guide](#)上，你可以查看到不同SIMD指令的CPI。注意到div类指令的CPI为8，远高于mul类指令的CPI 0.5，而编译器通常不会自动将SIMD的div指令优化为mul指令，因此我们需要手动将div 4改为mul 0.25。此外，还需要注意使用SIMD时，需要处理不能被4除尽的\$N\$的情况。这时你需要一些\$ISD\$的操作。
- 本项目无需Cache Blocking，原因见Profile结果（后文）。
- 最后，为每个可能加入OpenMP的循环引入OpenMP。
- 样例代码：

```
void impl(int N, int step, double *p) {
double divisor[4] = {
    0.25f,
    0.25f,
    0.25f,
    0.25f,
};
__m256d p_divisor = _mm256_loadu_pd(divisor);

// rearrange
int N2 = (N + 1) / 2;
double *p_part[2] = {
    aligned_alloc(32, N2 * N * sizeof(double)),
    aligned_alloc(32, N2 * N * sizeof(double)),
};
#pragma omp parallel for
for (int i = 0; i < N; ++i) {
    int part = i & 1;
    for (int j = 0; j < N; ++j) {
        p_part[part][i * N2 + j / 2] = p[i * N + j];
        part ^= 1;
    }
}

// caculate
int INPUTpartID = 1;
int OUTPUTpartID = 0;

if (N & 1) { // N = odd
    for (int k = 0; k < step; k++) {
#pragma omp parallel for
        for (int i = 1; i < N - 1; i++) {
            int j_head = (INPUTpartID + i) & 1;
            int j_begin = i * N2 + j_head;
            int j_end = (i + 1) * N2 - 1;
            int j = j_begin;
            for (; j < j_end - 3; j += 4) {
                __m256d p1 = _mm256_loadu_pd(&p_part[INPUTpartID][j - N2]);
                __m256d p2 = _mm256_loadu_pd(&p_part[INPUTpartID][j - j_head]);
                __m256d p3 = _mm256_loadu_pd(&p_part[INPUTpartID][1 + j - j_head]);
                __m256d p4 = _mm256_loadu_pd(&p_part[INPUTpartID][j + N2]);
                __m256d sum1 = _mm256_add_pd(p1, p2);
                __m256d sum2 = _mm256_add_pd(p3, p2);
                __m256d sum3 = _mm256_add_pd(sum1, sum2);
                __m256d result = _mm256_mul_pd(sum3, p_divisor);
                __mm256_storeu_pd(&p_part[OUTPUTpartID][j], result);
            }

            // for the tail
            for (; j < j_end; j++) {
                double p1 = p_part[INPUTpartID][j - N2];
                double p2 = p_part[INPUTpartID][j - j_head];
                double p3 = p_part[INPUTpartID][1 + j - j_head];
                double p4 = p_part[INPUTpartID][j + N2];
                p_part[OUTPUTpartID][j] = (p1 + p2 + p3 + p4) / 4.0f;
            }
        }

        int temp = INPUTpartID;
        INPUTpartID = OUTPUTpartID;
        OUTPUTpartID = temp;
    }
} else { // N = even
    for (int k = 0; k < step; k++) {
#pragma omp parallel for
        for (int i = 1; i < N - 1; i++) {
            int j_head = (INPUTpartID + i) & 1;
            int j_begin = i * N2 + j_head;
            int j_end = N2 - 1 + j_begin;
            int j = j_begin;
            for (; j < j_end - 3; j += 4) {
                __m256d p1 = _mm256_loadu_pd(&p_part[INPUTpartID][j - N2]);
                __m256d p2 = _mm256_loadu_pd(&p_part[INPUTpartID][j - j_head]);
                __m256d p3 = _mm256_loadu_pd(&p_part[INPUTpartID][1 + j - j_head]);
                __m256d p4 = _mm256_loadu_pd(&p_part[INPUTpartID][j + N2]);
                __m256d sum1 = _mm256_add_pd(p1, p2);
                __m256d sum2 = _mm256_add_pd(p3, p4);
                __m256d sum3 = _mm256_add_pd(sum1, sum2);
                __m256d result = _mm256_mul_pd(sum3, p_divisor);
                __mm256_storeu_pd(&p_part[OUTPUTpartID][j], result);
            }

            // for the tail
            for (; j < j_end; j++) {
                double p1 = p_part[INPUTpartID][j - N2];
                double p2 = p_part[INPUTpartID][j - j_head];
                double p3 = p_part[INPUTpartID][1 + j - j_head];
                double p4 = p_part[INPUTpartID][j + N2];
                p_part[OUTPUTpartID][j] = (p1 + p2 + p3 + p4) / 4.0f;
            }
        }

        int temp = INPUTpartID;
        INPUTpartID = OUTPUTpartID;
        OUTPUTpartID = temp;
    }
}

// rearrange back
#pragma omp parallel for
for (int i = 0; i < N; ++i) {
    int part = i & 1;
    for (int j = 0; j < N; ++j) {
        p[i * N + j] = p_part[part][i * N2 + j / 2];
        part ^= 1;
    }
}

free(p_part[0]);
free(p_part[1]);
}
```

附：Profile（一会更）

[Click here](#) to view. Search or link to this question with @514.

Sign up for more classes at <http://piazza.com/shanghaitech.edu.cn>.

Thanks,
The Piazza Team
--
Contact us at team@piazza.com

You're receiving this email because [wanghi2022@shanghaitech.edu.cn](#) is enrolled in CS 110 at ShanghaiTech University. Click [here](#) to unsubscribe from digest emails. Or, [sign in](#) to manage your email preferences or [un-enroll](#) from this class.