

# GPU 编程作业

学号：2120210492 姓名：肖飞

## 一、OneAPI tools 培训内容

2023  
英特尔® oneAPI  
校园黑客松竞赛

目录

基于Intel oneAPI编译器/编译选项的性能优化技巧

- 编译原理
- 编译器工作方式
- 常用优化编译选项
- 常用性能优化技巧

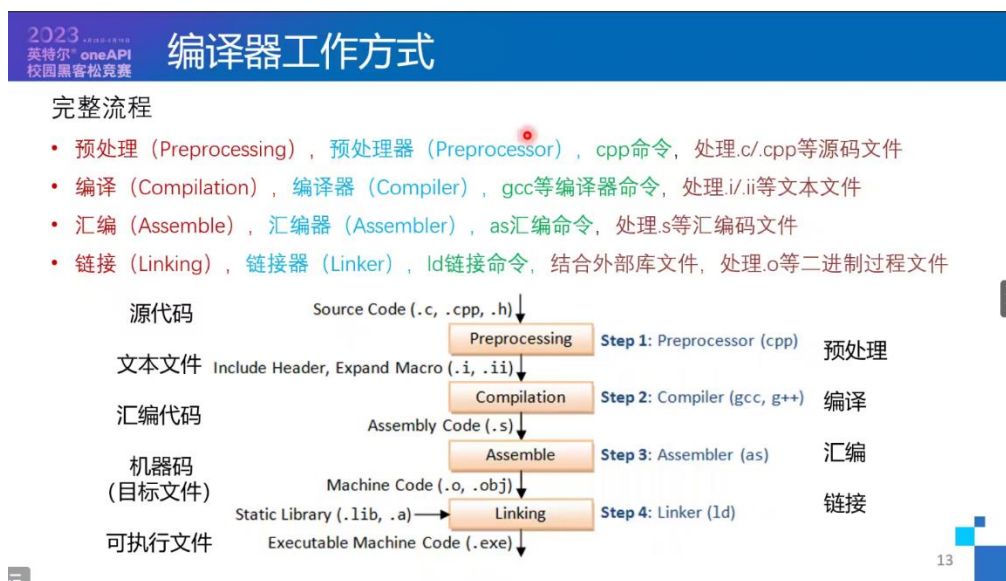
Intel oneAPI VTune/ITAC分析工具使用方法

- 编译原理
- 编译器工作方式
- 常用编译选项
- 常用性能优化技巧

三

5

整体内容主要分为两部分，第一部分是基于 Intel oneAPI 编译器/编译选项的性能优化技巧，第二部分是如何通过分析工具对程序进行性能分析，从而找到可以进行优化的地方。



首先简单介绍了编译的过程。

面向 HPC 的英特尔® oneAPI 工具

英特尔® oneAPI

HPC 工具套件

提供可扩展的快速应用

简介

英特尔® oneAPI 基础工具套件中的一个工具套件，用于在 C++、Fortran、OpenMP 和 MPI 上构建高性能、可扩展的并行代码，支持从企业到云以及从 HPC 到 AI 等各种应用。

这款产品的适用对象

- OEM/ISV
- C++、Fortran、OpenMP 和 MPI 开发人员

重要性

- 提升英特尔® 至强® 和酷睿® 处理器以及英特尔® 加速器的性能
- 更轻松交付基于行业标准的快速、可扩展、可靠的并行代码

英特尔® oneAPI 基础和 HPC 工具套件

直接编程

英特尔® C++ 编译器 (经典版)

英特尔® Fortran 编译器 (经典版)

英特尔® Fortran 编译器 (测试版)

英特尔® oneAPI DPC++/C++ 编译器

英特尔® DPC++ 兼容性工具

英特尔® Python 分发版

面向 oneAPI 基础工具套件的英特尔® FPGA 插件

基于 API 的编程

英特尔® MPI 函数库

英特尔® oneAPI DPC++ 库 - oneDPL

英特尔® oneAPI 数学核心函数库 - oneMKL

英特尔® oneAPI 数据解析库 - oneDAL

英特尔® oneAPI 线程构建模块 - oneTBB

英特尔® oneAPI 视频处理库 - oneVPL

英特尔® oneAPI 聚合通信库 - oneCCL

英特尔® oneAPI 深度学习网络库 - oneNN

英特尔® 生成性能基元 - 英特尔® IPP

分析和调试工具

英特尔® Inspector

英特尔® 跟踪分析器和采生器

英特尔® 集群检查器

英特尔® VTune™ 分析器

英特尔® Advisor

英特尔® GDB 分发版

英特尔® oneAPI HPC 工具套件 + 英特尔® oneAPI 基础工具套件

intel oneAPI HPC TOOLKIT

了解更多信息: [intel.com/oneAPI-HPCKit](https://intel.com/oneAPI-HPCKit)

面向高性能计算工作负载的特定领域工具套件 intel

Intel 的工具套件。

2023 英特尔® oneAPI 校园黑客松竞赛

常用优化编译选项

- 自动编译优化级别选项 -O
  - ✓ -O0  
不进行优化
  - ✓ -O1  
启用速度优化并禁用一些会增加代码大小并影响速度的优化。  
禁用某些内部函数的内联。对于代码非常大、分支很多、执行时间不受循环内代码支配的应用程序，O1 选项可以提高性能。
  - ✓ -O2  
默认设置，启用速度优化。  
向量化在 O2 及更高级别启用。启用函数内联，推测循环展开，部分冗余消除，及针对代码速度进行优化的选项。
  - ✓ -O3  
执行 O2 优化并开启更激进优化策略。依赖分析可能会导致更长的编译时间。O3 优化可能不会带来更高的性能，在某些情况下可能会减慢代码的速度。对于具有大量使用浮点计算和处理大型数据集的循环的应用程序，建议使用 O3 选项。

2023 英特尔® oneAPI 校园黑客松竞赛

常用优化编译选项

- 浮点计算控制 -fp-model
  - ✓ precise  
关闭浮点数值不安全的优化
  - ✓ fast[=1|2]  
启用更激进的浮点数优化策略
  - ✓ source  
将中间结果与源数据保持对齐
- zmm 寄存器控制 -qopt-zmm-usage=
  - ✓ low  
指定编译器尽量少的使用 zmm 寄存器
  - ✓ high  
指定编译器不对 zmm 寄存器限制使用

自动编译时进行一些优化选项的设置可能在不修改代码的情况下提高程序性能。

2023  
英特尔® oneAPI  
校园黑客松竞赛

常用性能优化技巧

- 从编程方式出发
  - ✓ 访存/寄存器优化
    - 调整数据访问顺序，利用局部性原理。  
C语言数据行优先存储、Fortran语言数据列优先存储.....
    - Cache line数据对齐
    - 指令预取prefetch
  - ✓ 冗余优化
    - 分析算法功能，减少不必要代码，或删除或合并
  - ✓ 算法优化
    - 利用更好性能的函数实现相同/类似效果的功能  
Intel® oneMKL、oneDNN.....
- 从编译器/编译选项出发
  - ✓ 添加编译优化选项
  - ✓ 使用Intel®编译器

Row-major order

Column-major order

16

第二部分介绍的分析工具主要有以下几个：

Intel 分析工具

Intel® Advisor

Intel® VTune™ Profiler

Intel® Inspector

Intel® Trace Analyzer & Collector

Intel® Cluster Checker

## 用于GPU计算分析的Intel分析工具

### Intel® Advisor

高效矢量化，线程，内存使用和加速器分流的设计代码

#### Offload Advisor

- 找出有价值的分流机会
- 检测瓶颈和关键边界因素
- 通过对性能，空间和瓶颈进行建模，甚至在拥有硬件之前就可以准备好代码

#### 屋顶线分析

- 查看针对硬件限制的性能余量
- 通过确定瓶颈以及哪些优化将带来最大收益来确定性能优化策略
- 可视化优化进度

#### 流程图分析器

- 可视化CPU / GPU代码并获取有关CPU设备的建议

### Intel® VTune™ Profiler

快速查找并修复性能瓶颈，实现硬件的所有价值

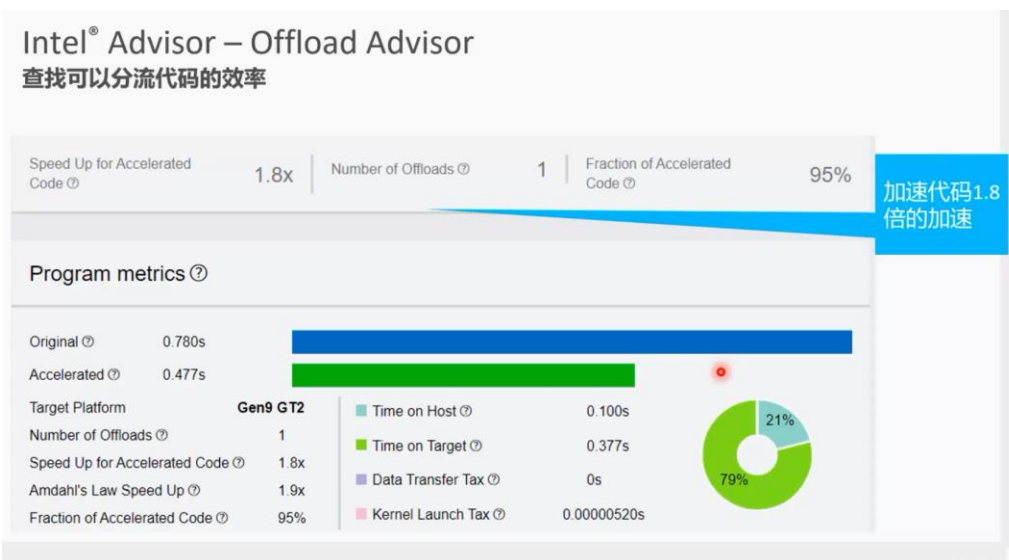
#### 分流性能调试

- 探索平台上各种CPU和GPU内核上的代码执行
- 关联CPU和GPU活动
- 识别应用程序是否受GPU或CPU限制

#### GPU计算热点

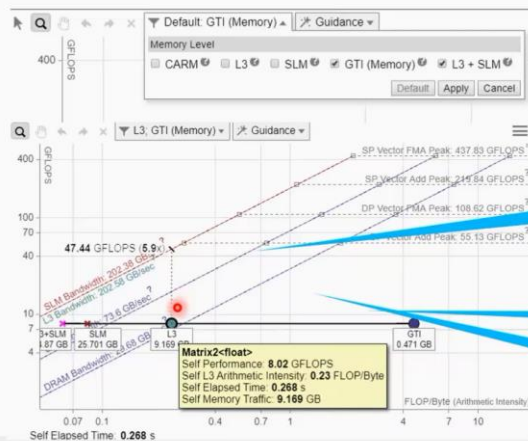
- 分析最耗时的GPU内核，基于GPU硬件指标表征GPU使用情况
- 源代码行级别和内核汇编级别的GPU代码性能

Advisor 对于 python 的性能优化帮助可能不大，而 VTune 对于 python 代码也可以进行分析。



程序在 GPU 上的加速效果

## Intel® Advisor – GPU Roofline 寻找有效的优化策略



配置要显示的级别

显示每个循环的性能余量

可能的瓶颈

建议下一步优化

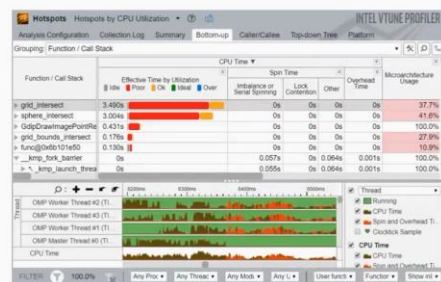
## 分析和调整应用程序性能 Intel® VTune™ Profiler

### 节省优化代码的时间

- 准确地分析C, C++, Fortran, Python, Go, Java或任何组合
- 优化CPU, 线程, 内存, 缓存, 存储等
- 节省时间: 丰富的分析可带来洞察力

### 2021.1 版本新功能 (部分列表)

- Platform Profiler产品发布 - 收集更多的指标
- 设计和优化Intel® Optane™ DC持久性存储器
- 应用性能快照 - 添加了通信模式诊断, 分析更多进程
- Linux - Perf的广泛使用使分析无需添加驱动程序



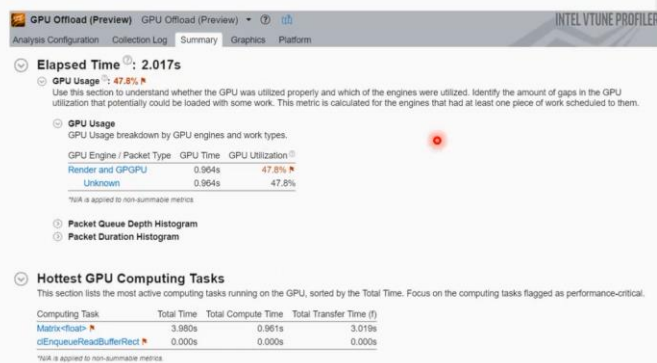
了解更多: [software.intel.com/vtune](https://software.intel.com/vtune)

可以分析出程序中不同函数的时间等开销及可优化的空间。

## Intel® VTune™ Profiler – GPU Offload 优化GPU使用率

### 可以分析

- 确定应用程序使用DPC++或OpenCL内核的效率
- 分析一段时间内Intel® Media SDK任务的执行情况 (仅适用于Linux目标)
- 探索GPU的使用情况并在每个时刻分析GPU引擎的软件队列

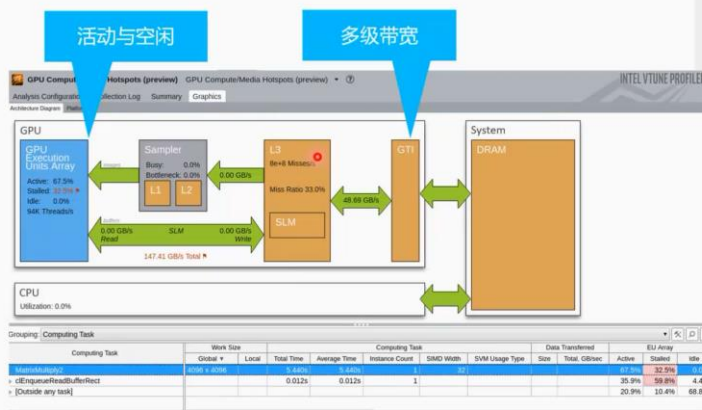




## Intel® VTune™ Profiler – GPU热点 优化GPU使用率

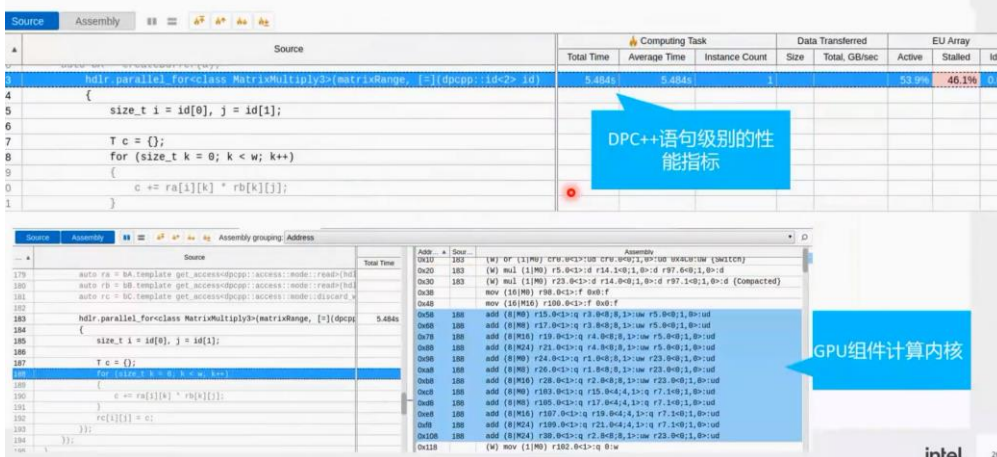
运行VTune Profiler GPU热点以识别GPU使用率低和停滞的原因

选择“图形”选项卡以查看体系结构的高级图



## GPU 分析

## Intel® VTune™ DPC++ 代码分析



## 具体代码语句的性能分析

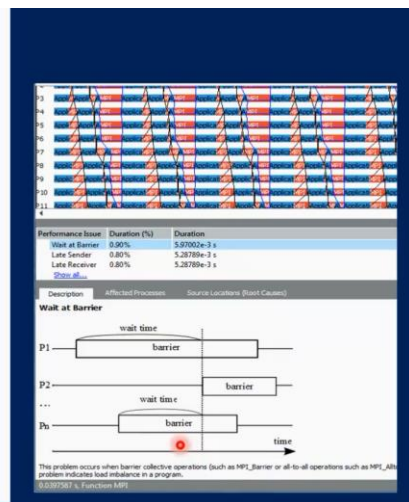
## Intel® Trace Analyzer & Collector 分析并可视化分析MPI应用程序行为

### 扩展MPI应用

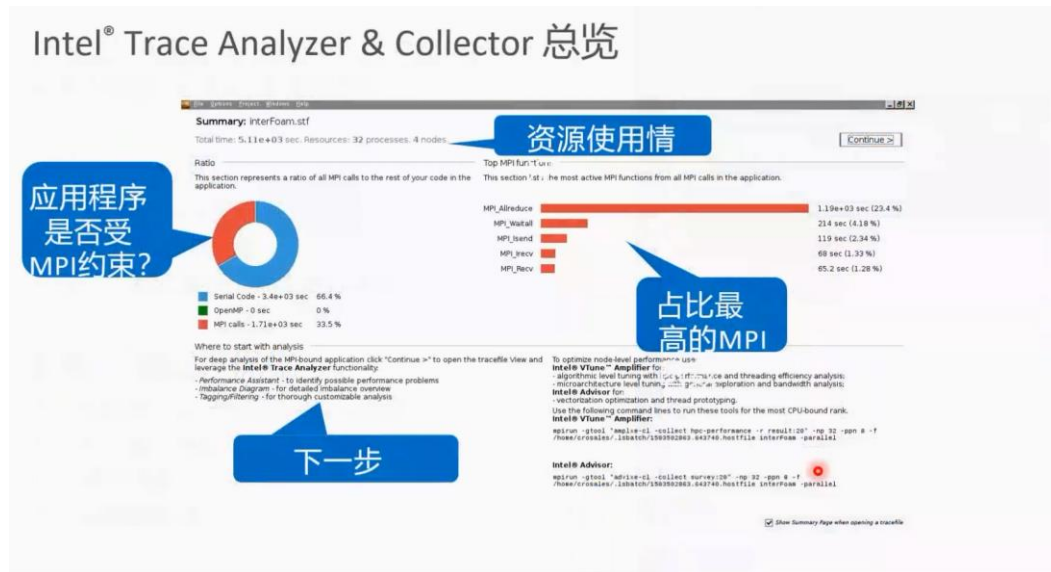
- 扩展性能: 在更多节点上运行
- 向前扩展: 多核就绪
- 有效扩展: 在更多节点上进行调试

### 分析, 调试, 优化

- 可视化并了解并行应用程序的行为
- 评估性能分析统计信息和负载平衡
- 分析常见的MPI问题
- 识别通信热点



## MPI 程序行为分析



## 二、动手练习

实验环境如下：

实验平台：x64 平台

CPU：12th Gen Intel(R) Core(TM) i5-12500H 3.10 GHz

软件：Intel® oneAPI Base Toolkit & HPC Toolkit

操作系统：Windows 11 系统

编译器：TDM-GCC & ICC 编译器

下面是一段矩阵加法的 C 代码：

```

#include <stdio.h>
#include <time.h>
#include <stdlib.h>

#define N 20000
#define FIRST i
#define SECOND j

int main()
{
    int i,j,k;
    int **a,**b,**c;
    clock_t malloc_start,malloc_end,init_start,init_end,compute_start,compute_end;
    double malloc_time,init_time,compute_time;
    malloc_start=clock();
    printf("Starting malloc!\n");
    a=(int **)malloc(sizeof(int *) * N);
    b=(int **)malloc(sizeof(int *) * N);
    c=(int **)malloc(sizeof(int *) * N);
    printf("Starting malloc for !\n");
    for(i=0;i<N;i++)
    {
        a[i]=(int *)malloc(sizeof(int) * N);
        b[i]=(int *)malloc(sizeof(int) * N);
        c[i]=(int *)malloc(sizeof(int) * N);
    }
    malloc_end=clock();
    init_start=clock();
    printf("Starting init ! \n");
    for(FIRST=0;FIRST<N;FIRST++)
    {
        "array-ij.c" 64L, 1645C

```

```

        malloc_end=clock();
        init_start=clock();
        printf("Starting init ! \n");
        for (FIRST=0;FIRST<N;FIRST++)
        {
            for (SECOND=0;SECOND<N;SECOND++)
            {
                a[i][j]=i;
                b[i][j]=j;
                c[i][j]=0;
            }
        }
        init_end=clock();
        printf("Starting compute ! \n");
        compute_start=clock();
        for (FIRST=0;FIRST<N;FIRST++)
        {
            for (SECOND=0;SECOND<N;SECOND++)
            {
                c[i][j]=a[i][j] + b[i][j];
            }
        }
        compute_end=clock();
        malloc_time=(double)(malloc_end-malloc_start) / CLOCKS_PER_SEC;
        init_time=(double)(init_end-init_start) / CLOCKS_PER_SEC;

```



```

compute_end=clock();
malloc_time=(double)(malloc_end-malloc_start) / CLOCKS_PER_SEC;
init_time=(double)(init_end-init_start) / CLOCKS_PER_SEC;
compute_time=(double)(compute_end-compute_start) / CLOCKS_PER_SEC;
printf("malloc time is %lf,init time is %lf,compute time is %lf\n",malloc_time,init_time,compute_time);
for(i=0;i<N;i++)
{
    free(a[i]);
    free(b[i]);
    free(c[i]);
}
free(a);
free(b);
free(c);
return 0;

```

按照不同策略对上述代码进行优化，并编译运行测试不同阶段的时间开销，得到以下结果。

```

[sch0205@ln2%bscc-t6 code]$ vim array-
array-gcc-ij.exe    array-gcc-ij-03.exe  array-gcc-ji-02.exe  array-icc-ij.exe    array-icc-ji.exe    array-ij.c
array-gcc-ij-02.exe array-gcc-ji.exe     array-gcc-ji-03.exe  array-icc-ij-03.exe array-icc-ji-03.exe array-ji.c
[sch0205@ln2%bscc-t6 code]$ vim array-ij.c
[sch0205@ln2%bscc-t6 code]$ ./array-gcc-ij.exe
Starting malloc!
Starting malloc for !
Starting init !
Starting compute !
malloc time is 0.140000,init time is 3.850000,compute time is 4.920000
[sch0205@ln2%bscc-t6 code]$ vim array-ji.c
[sch0205@ln2%bscc-t6 code]$ ./array-gcc-ji.exe
Starting malloc!
Starting malloc for !
Starting init !
Starting compute !
malloc time is 0.150000,init time is 24.310000,compute time is 20.470000
[sch0205@ln2%bscc-t6 code]$ vim array-ij.c
[sch0205@ln2%bscc-t6 code]$

```

对于 C 程序，遵循行优先访存（第一次）的性能相比按照列优先访存（第二次）的性能由有较大的提升。

```

[sch0205@ln2%bscc-t6 code]$ vim array-
array-gcc-ij.exe    array-gcc-ij-03.exe  array-gcc-ji-02.exe  array-icc-ij.exe    array-icc-ji.exe    array-ij.c
array-gcc-ij-02.exe array-gcc-ji.exe     array-gcc-ji-03.exe  array-icc-ij-03.exe array-icc-ji-03.exe array-ji.c
[sch0205@ln2%bscc-t6 code]$ vim array-ij.c
[sch0205@ln2%bscc-t6 code]$ ./array-gcc-ij.exe
Starting malloc!
Starting malloc for !
Starting init !
Starting compute !
malloc time is 0.140000,init time is 3.850000,compute time is 4.920000
[sch0205@ln2%bscc-t6 code]$ vim array-ji.c
[sch0205@ln2%bscc-t6 code]$ ./array-gcc-ji.exe
Starting malloc!
Starting malloc for !
Starting init !
Starting compute !
malloc time is 0.150000,init time is 24.310000,compute time is 20.470000
[sch0205@ln2%bscc-t6 code]$ vim array-ij.c
[sch0205@ln2%bscc-t6 code]$ ./array-gcc-ij-02.exe
Starting malloc!
Starting malloc for !
Starting init !
Starting compute !
malloc time is 0.210000,init time is 2.910000,compute time is 0.810000
[sch0205@ln2%bscc-t6 code]$

```

使用了 -O2 编译优化选项后，各阶段代码性能也有了一定提升。

```

[sch0205@ln2%bscc-t6 code]$ vim array-ji.c
[sch0205@ln2%bscc-t6 code]$ ./array-gcc-ji.exe
Starting malloc!
Starting malloc for !
Starting init !
Starting compute !
malloc time is 0.150000,init time is 24.310000,compute time is 20.470000
[sch0205@ln2%bscc-t6 code]$ vim array-ij.c
[sch0205@ln2%bscc-t6 code]$ ./array-gcc-ij-O2.exe
Starting malloc!
Starting malloc for !
Starting init !
Starting compute !
malloc time is 0.210000,init time is 2.910000,compute time is 0.810000
[sch0205@ln2%bscc-t6 code]$ ./array-gcc-ji-O2.exe
Starting malloc!
Starting malloc for !
Starting init !
Starting compute !
malloc time is 0.160000,init time is 32.990000,compute time is 20.920000
[sch0205@ln2%bscc-t6 code]$

```

对于列优先访存的代码，使用-O2 选项编译反倒增加了时间开销。

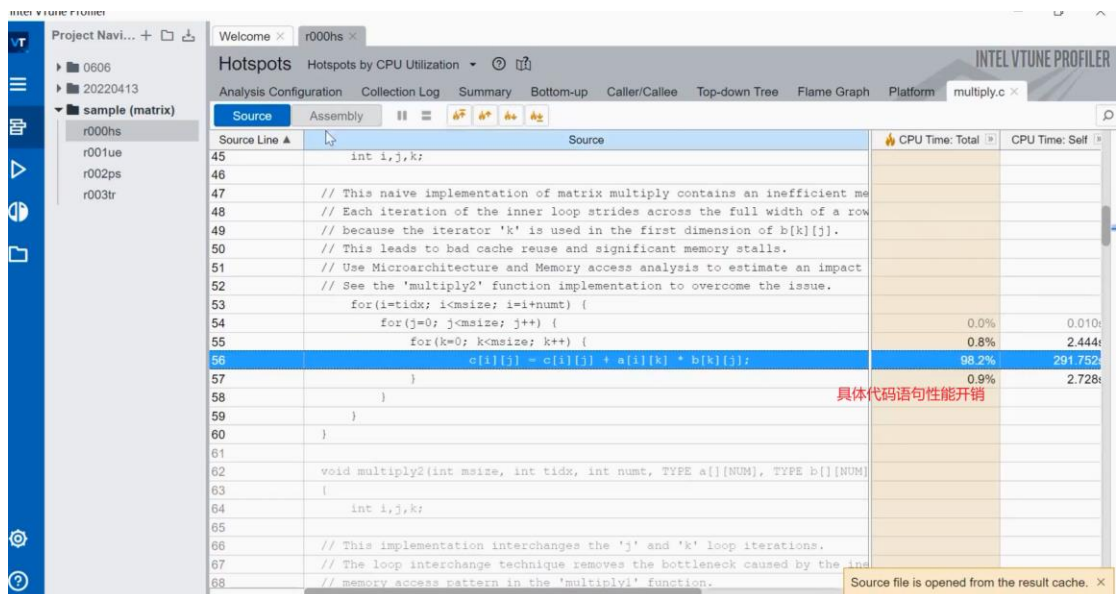
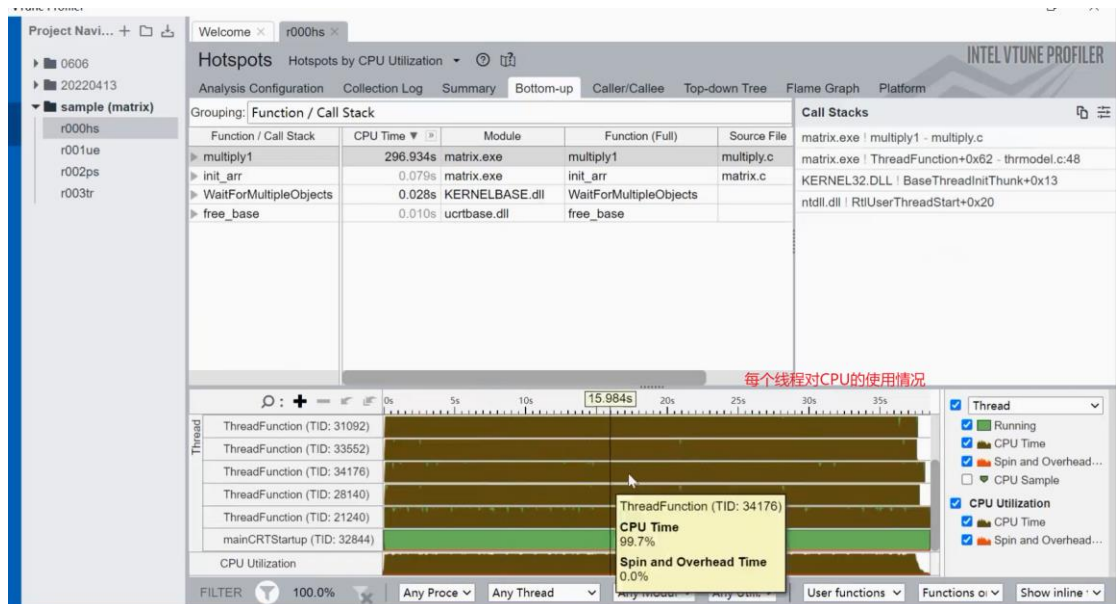
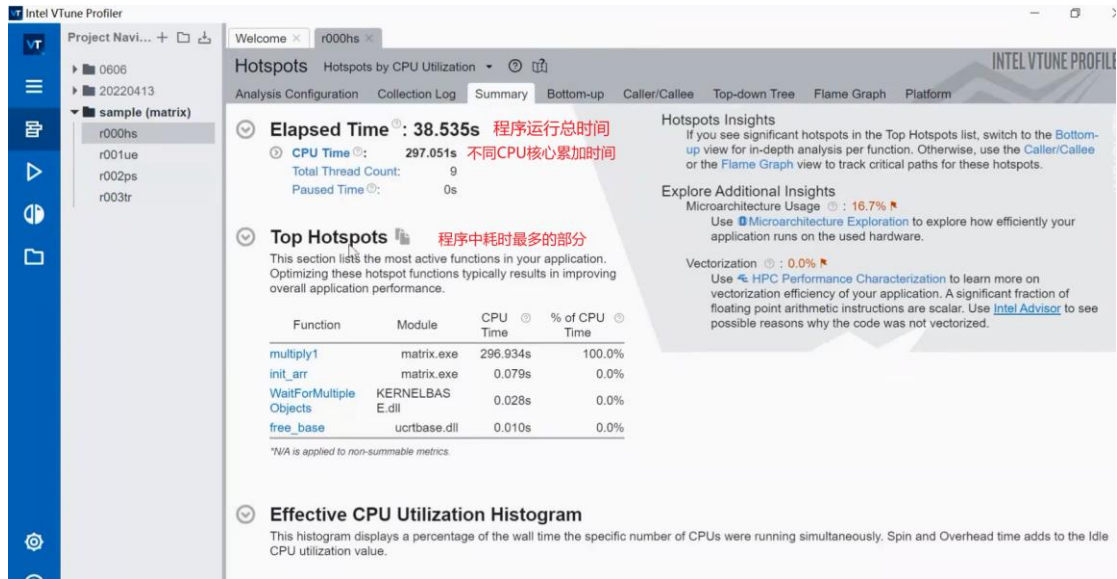
```

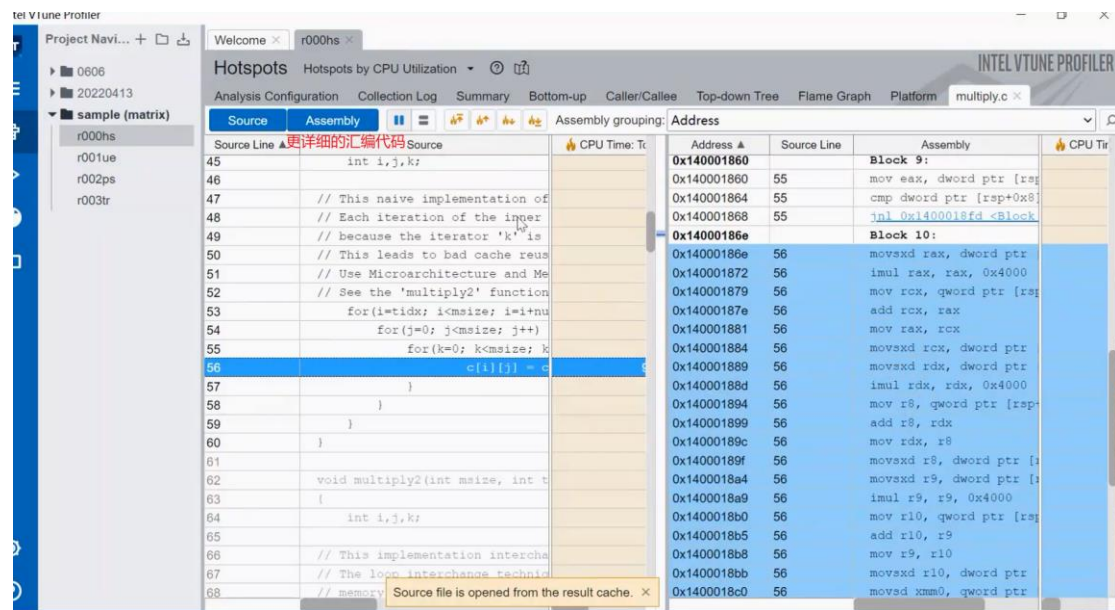
[sch0205@ln2%bscc-t6 code]$ ls
array-gcc-ij.exe      array-gcc-ji.exe      array-icc-ij.exe      array-icc-ji-O3.exe  compile.sh  run.sh
array-gcc-ij-O2.exe   array-gcc-ji-O2.exe   array-icc-ij-O3.exe   array-ij.c           matrix-c
array-gcc-ij-O3.exe   array-gcc-ji-O3.exe   array-icc-ji.exe      array-ji.c           matrix.c
[sch0205@ln2%bscc-t6 code]$ ./array-icc-ij.exe
Starting malloc!
Starting malloc for !
Starting init !
Starting compute !
malloc time is 0.150000,init time is 2.670000,compute time is 0.950000
[sch0205@ln2%bscc-t6 code]$ ./array-icc-ji.exe
Starting malloc!
Starting malloc for !
Starting init !
Starting compute !
malloc time is 0.150000,init time is 3.380000,compute time is 1.820000
[sch0205@ln2%bscc-t6 code]$ ./array-icc-ij-O3.exe
Starting malloc!
Starting malloc for !
Starting init !
Starting compute !
malloc time is 0.160000,init time is 2.830000,compute time is 0.930000
[sch0205@ln2%bscc-t6 code]$ ./array-icc-ji-O3.exe
Starting malloc!
Starting malloc for !
Starting init !
Starting compute !
malloc time is 0.150000,init time is 3.030000,compute time is 0.840000
[sch0205@ln2%bscc-t6 code]$

```

将 gcc 编译器更换为 Intel 的 icc 编译器后，同样的代码性能也有了一定提升，即便将行优先访存更改为列优先访存，性能也只是下降了一点，也就是说，使用 Intel 编译器可以在考虑较少代码优化细节（如访存顺序、编译优化选项等）的情况下仍旧保持不错的性能。

使用 VTune 工具（windows 版本）进行代码性能分析：





### 三、项目地址

Git 项目链接: [https://github.com/XFLasdf/Guass\\_GPU](https://github.com/XFLasdf/Guass_GPU)