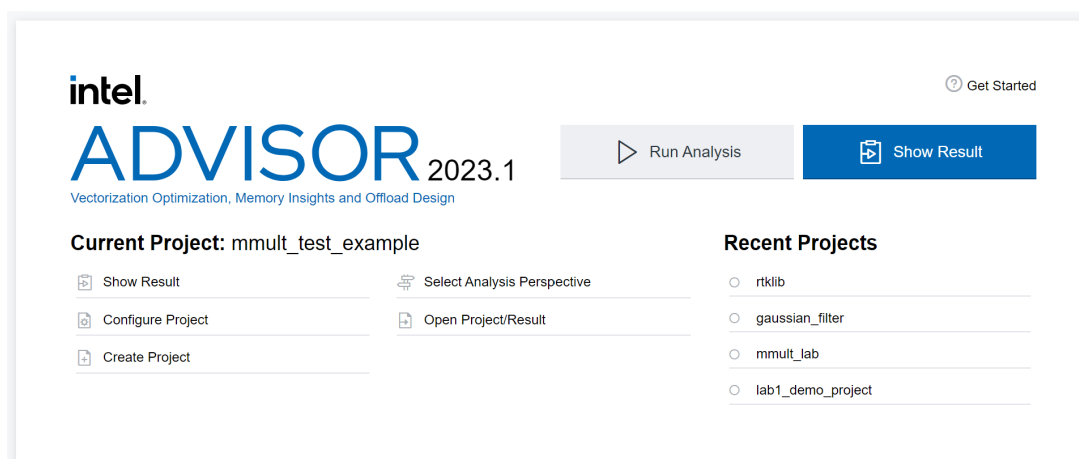


## Lab2. Intel advisor

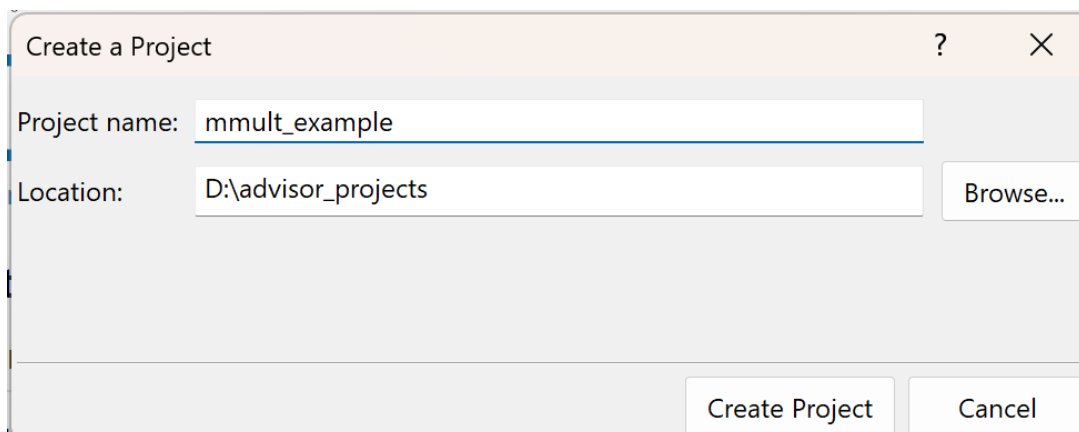
Этап с настройкой окружения для запуска Intel Advisor соответствует аналогичному шагу в предыдущей лабораторной работе. Запуск Intel Advisor с использованием пользовательского интерфейса (выполняется из консоли с настроенным окружением):

```
C:\Users\k.sandalov>advisor-gui
```

Для создания нового проекта - **Create project**:



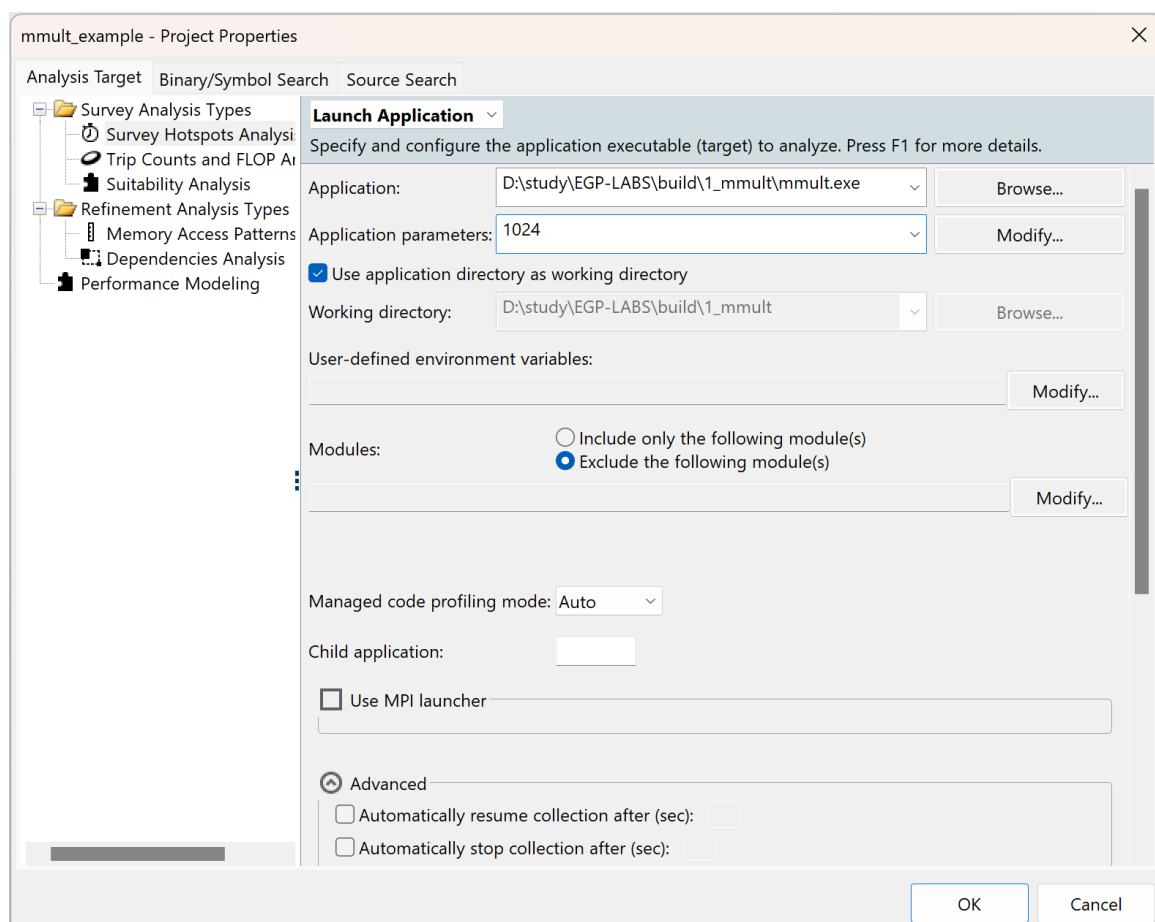
Выбор имени проекта и его расположение:



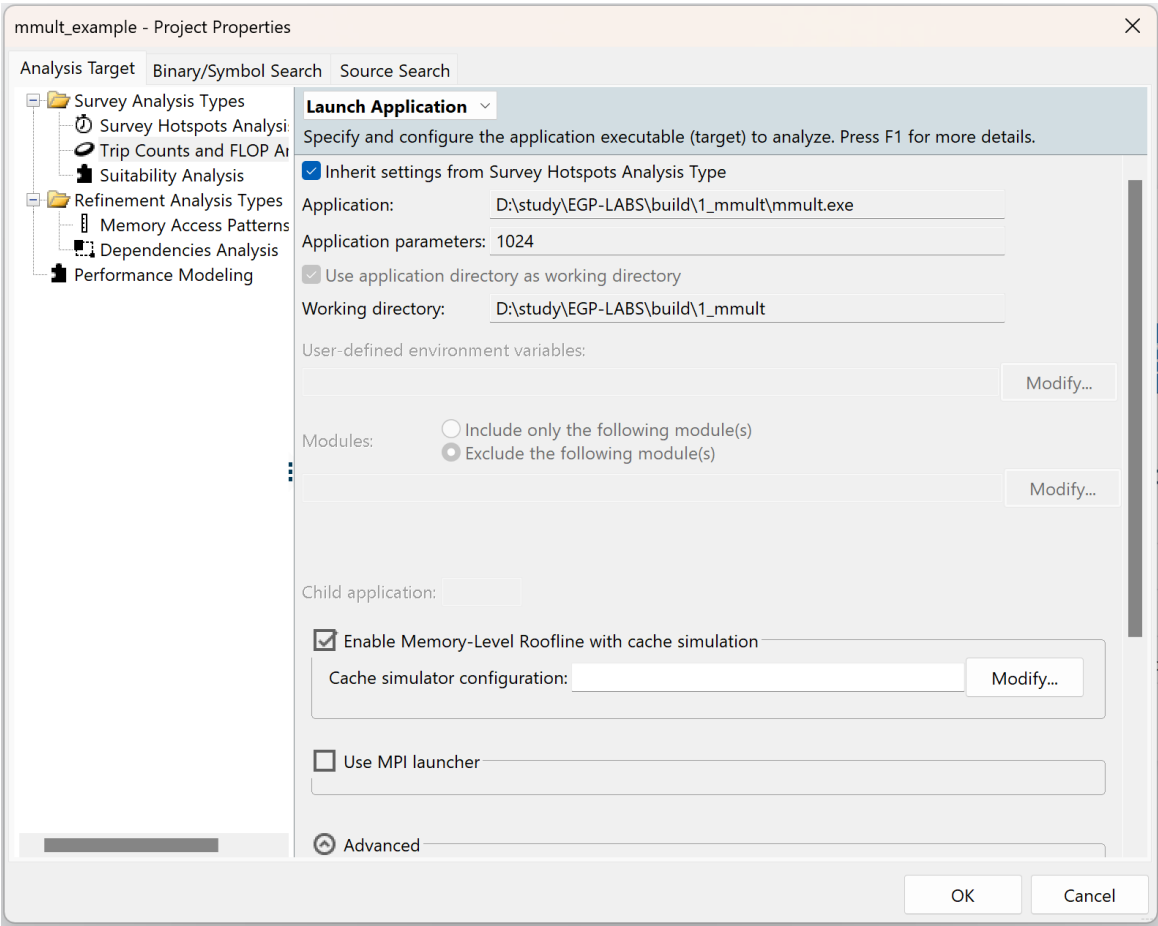
После этого, для начальной настройки проекта нужно на вкладке **Analysis Target** -> **Survey Hotspots Analysis** выбрать анализируемое приложение, указав полный путь до полученного .exe файла в окне **Application**. Соответственно, аргументы запускаемого приложения (выбран режим **Launch Application**), задаются в окне **Application parameters**.

В рамках лабораторной работы используем тестовое приложение mmult из этого репозитория. Сперва необходимо собрать его базовую версию со следующими опциями компиляции:

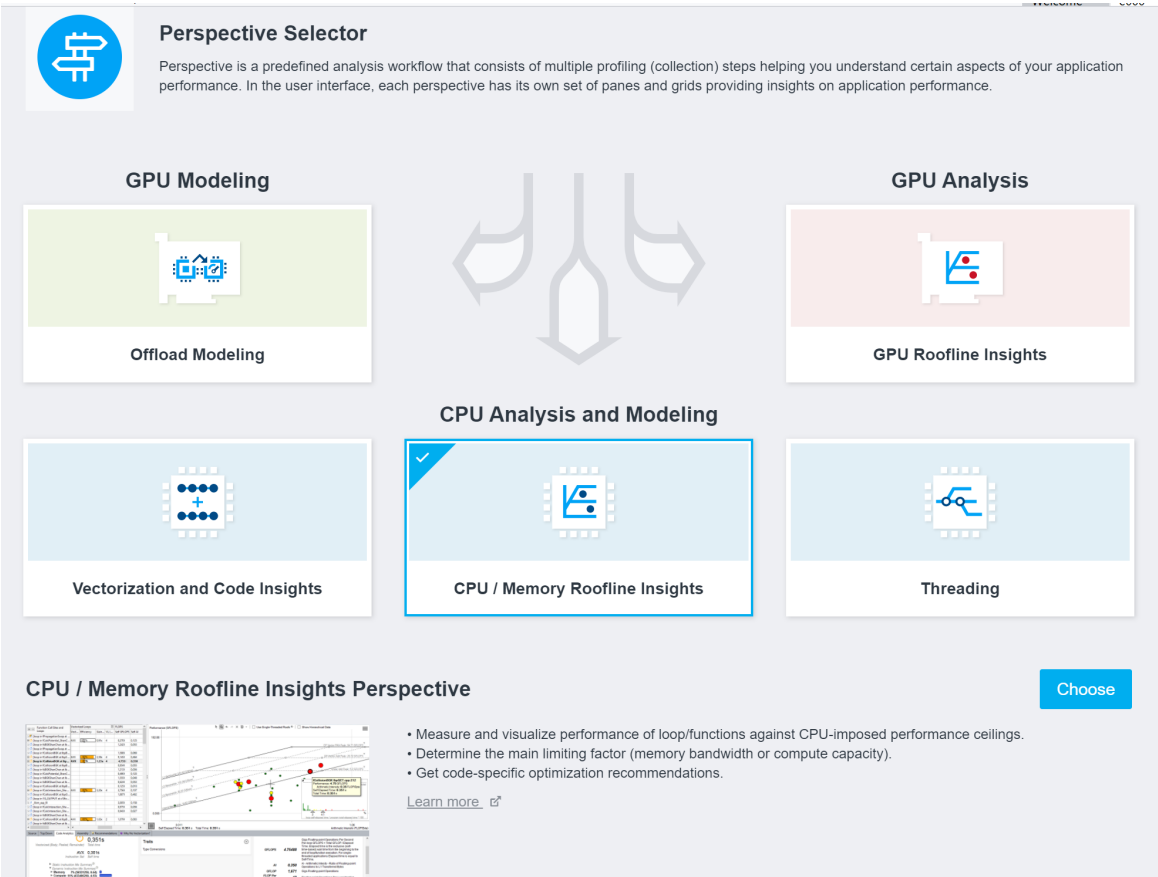
```
target_compile_options(mmult PRIVATE
/I$ENV{ONEAPI_ROOT}\\advisor\\latest\\include)
```



Дополнительно на вкладке **Analysis Target** -> **Trip Counts and FLOP Analysis** необходимо выбрать пункт **Enable Memory-Level Roofline with cache simulation**:



При создании нового проекта первым должно появиться окно с выбором определенного типа анализа или же, **Perspective selector**:



Необходимо небольшое отступление, чтобы в общих чертах описать цель и средства анализа, выполняемого в Intel Advisor.

В соответствии со своим названием, Intel Advisor способен проанализировать приложение (исходный код, отлабочная информация + динамический анализ исполняемого приложения) и дать советы по его оптимизации.

Для анализа приложения в нем представлены различные flow, так называемые, перспективы, которые помогут раскрыть различные аспекты эффективности утилизации ресурсов вашей вычислительной платформы анализируемым приложением.

Соответственно, разные flow содержат различные наборы анализов из тех, что вообще представлены в инструменте:

- Survey

Необходим для первичного анализа (самый легковесный) горячих участков приложения и эффективности векторизации. Далее можно уже применять какие-то предложения по повышению эффективности векторизации, и обратить свое внимание на то, какие у вас есть хотспоты и попытаться их оптимизировать

- Characterization

Служит для дальнейшего определения возможностей оптимизации хотспотов. Он позволит собрать более подробную информацию о том, как именно выполняются найденные хотспоты. Если это какие-то циклы, то, можно определить количество реально выполненных итераций, для функций же можно получить информацию о количестве их вызовов (это все Trip Counts). Тут же можно узнать, сколько было выполнено Floating Point / Int операций (FLOP) и уже на основании двух этих анализов - Trip Counts и FLOP построить очень полезный график [Roofline](#). Он подскажет, чем реально ограничен тот или иной хотспот - подсистемой памяти, либо вычислительной пропускной способностью (это применение векторизации и параллелизма), и отталкиваясь от этого уже можно будет думать, а нужно ли как-то улучшать чтение/запись данных или уже пытаться добавлять параллелизм, бороться с векторизацией.

- Memory Access Patterns

Используется для определения эффективности паттернов работы с памятью, ведь оптимальнее всего обращаться к памяти последовательно, либо же вообще попытаться избежать лишних операций чтения/записи в более дорогую память. Необходимо обозначить интересующие вас циклы, чтобы проанализировать, как в них происходит работа с памятью

- Dependencies

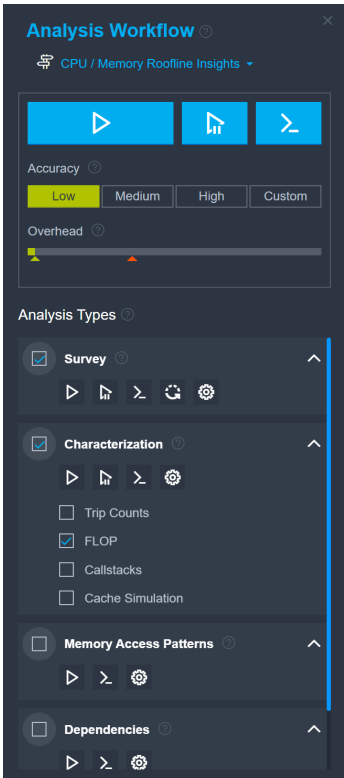
Если вы хотите добавлять параллелизм, то необходимо помнить, что эффективная параллельная обработка данных предполагает определенную (лучше всего конечно высокую) степень независимости между данными, обрабатываемыми в разных потоках. Этот анализ как раз таки позволяет определить, можно ли распараллелить тот или иной цикл по его итерациям эффективно, в зависимости от того, есть ли между его итерациями зависимости.

- Suitability

Когда мы поняли, что на предыдущих анализах у нас все отлично, мы максимально оптимально работаем с памятью и независимы по итерациям, можно попытаться добавить параллелизм. Однако тут тоже не все так просто и полезно было бы оценить, как будет масштабироваться наша параллельная программа с ростом кол-ва выполняемой работы и ее распределением по отдельным создаваемым потокам. Для этого необходимо расставить аннотации в соответствии с [рекомендациями](#).

Возвращаясь к разбору Intel Advisor, рассмотрим перспективу **CPU / Memory Roofline Insights**.

В блоке **Analysis Workflow** можно переключиться между разными перспективами (были представлены на предыдущем изображении) а так же подобрать параметры для текущего типа анализа и запустить сбор данных + непосредственно сам анализ.



Для первого запуска должен быть выбран **Survey + Characterization (FLOP)**, сравнить время выполнения этих анализов между собой.

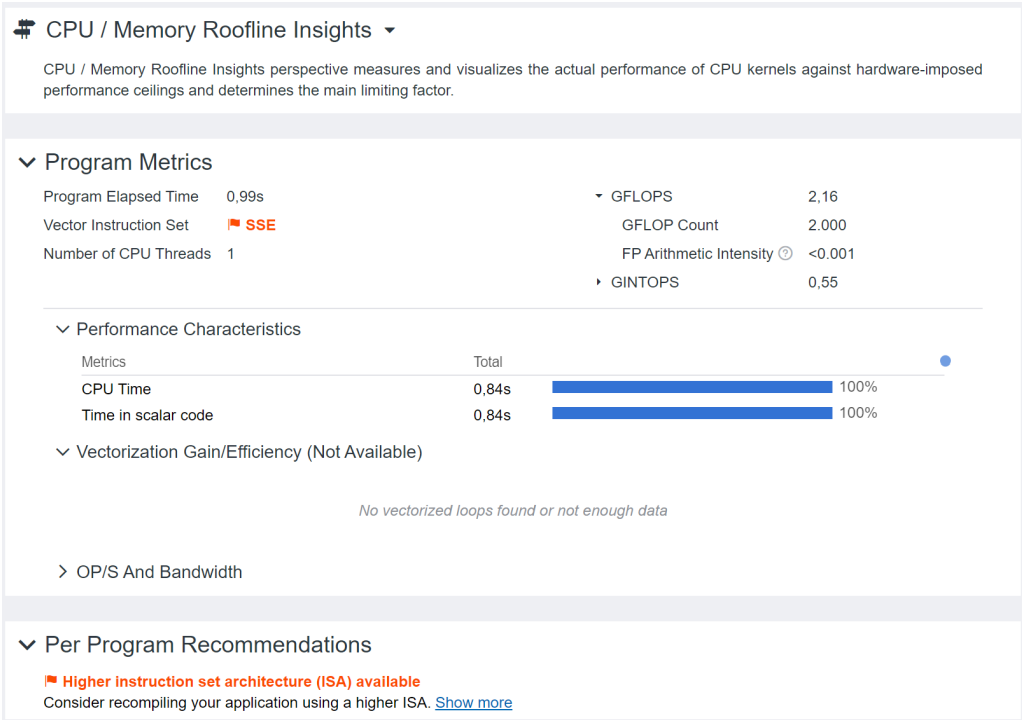
Задание: Результаты сравнения добавить в отчет. Время выполнения анализа можно найти на вкладке **Summary** в блоке **Collection Details**:

Collection Details	
Survey	
Collection started	21 February 2024, 03:06:16
Collection finished	21 February 2024, 03:06:17
Collection time	00 min 01 sec
Finalization time	00 min 01 sec
Full time	00 min 02 sec
Collection Log	<a href="#">See log</a>
Application Output	<a href="#">See output</a>
Collection Command Line	<a href="#">See command line</a>

**Summary:**

На этой вкладке приведены общие сведения о производительности приложения:

- Количество выполняемых потоков
- Время работы программы от старта выполнения первого потока до завершения последнего
- Используемый набор векторных инструкций
- Доля векторизованных вычислений
- GFLOPS = FLOPs / Seconds (аналогично и INTOPS)
- ArithmeticIntensity(AI) = FLOPs / Bytes (отношение кол-ва вычислений к кол-ву запрошенных для этого байтов из памяти)



- Теоретические пиковые значения производительности для данной платформы и реально задействованная мощность:

OP/S And Bandwidth

Effective OP/S And Bandwidth		Utilization	Hardware Peak
GFLOPS	<b>4.883</b>	7.76%	out of 62.940 (DP) GFLOPS
		4.00%	out of 121.938 (SP) GFLOPS
GINTOPS	<b>1.247</b>	2.88%	out of 43.286 (Int64) GINTOPS
		1.31%	out of 95.035 (Int32) GINTOPS
CPU <-> Memory [L1+NTS GB/s]	<b>29.326</b>	7.43%	out of 394.856 GB/s [bytes]

Тут же можно найти рекомендации, применимые глобально ко всему приложению, например: Использовать более широкий набор инструкций.

Задание: Зафиксировать время выполнения программы, факт использования векторных инструкций, текущие GFLOPS и AI для приложения. Пиковые значения пропускной способности вычислительной системы.

Топ самых тяжелых хотспотов и рекомендации для них. Подробнее про анализ хотспотов на вкладке **Survey and Roofline**:

Per Program Recommendations

Higher instruction set architecture (ISA) available

Consider recompiling your application using a higher ISA. [Show more](#)

Top Time-Consuming Loops

Loop	Self Time	Total Time
loop in multiply_simple<float> at mmult.cpp:45	0.820s	0.820s
loop in operator new at new_scalar.cpp:35	<0.001s	0.003s
loop in multiply_simple<float> at mmult.cpp:40	<0.001s	0.820s
loop in multiply_simple<float> at mmult.cpp:42	<0.001s	0.820s
loop in mmult_task at mmult.cpp:131	<0.001s	0.003s

Suitability And Dependencies Analysis Data

No data available. Collect Suitability or Dependency to see the results.

Recommendations

Target the AVX2 ISA

loop in multiply\_simple<float> at mmult.cpp:45

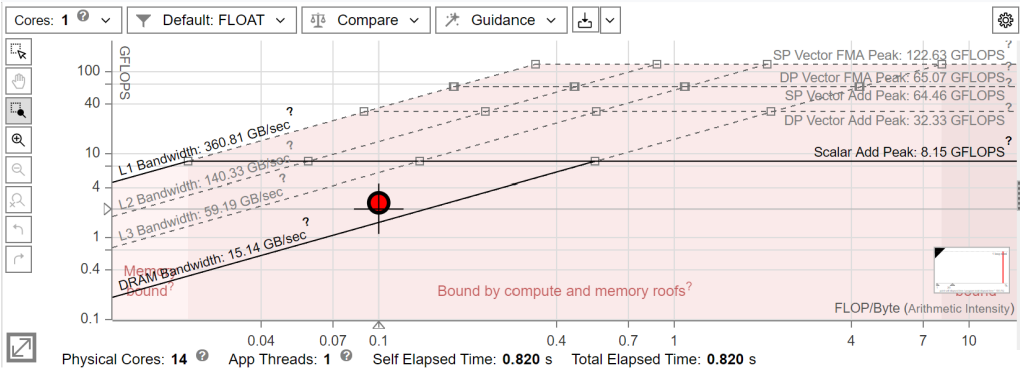
Remove system function call(s) inside loop

loop in operator new at new\_scalar.cpp:35

Remove system function call(s) inside loop

loop in mmult\_task at mmult.cpp:200

Roofline:



На графике руфлайна наглядно отображены ограничения пропускной способности (вычислительной - горизонтальные линии и подсистемы памяти - наклонные), а также то, насколько эффективно их использует анализируемое приложение. На графике отображаются точки, характеризующие производительность самых горячих участков тестового приложения. И в данном случае, это цикл на строке 45 в mmult.cpp. Для этой точки вычисляются значения GFLOPS/GINTOPS и AI, а так же, исходя из используемого набора инструкций, типа операций (FP или INT) и конфигурации подсистемы памяти ограничивающие ее крыши.

Соответственно, ограничением производительности будет служить  $\text{MIN}(\text{peakMemBandwidth} \times \text{AI}, \text{peakGFLOPS})$

Survey:

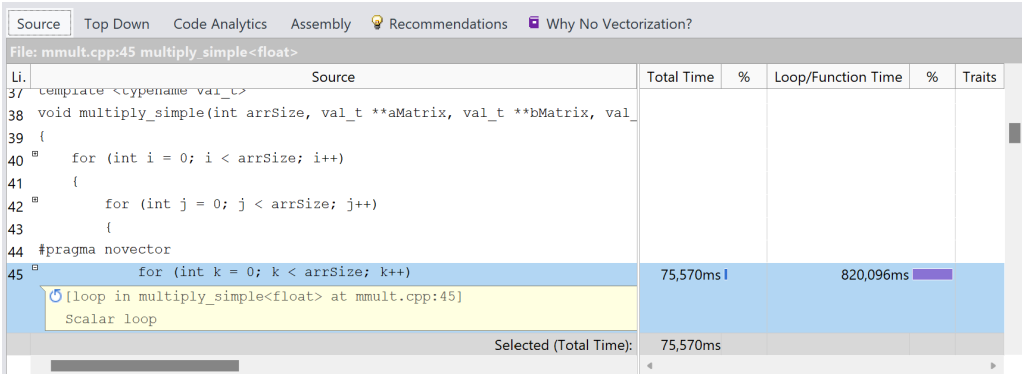
Function Call Sites and Loops	Performance Issues	CPU Time		Type	Why No Ve
		Total Time	Self Time		
<input checked="" type="checkbox"/> [loop in multiply_simple<float> at mmult.cpp:45]	<input checked="" type="checkbox"/> 1 Potential un...	0,820s	0,820s	Scalar	
<input checked="" type="checkbox"/> [loop in operator new at new_scalar.cpp:35]	<input checked="" type="checkbox"/> 1 System functi...	0,003s	0,000s	Scalar	
<input checked="" type="checkbox"/> mainCRTStartup	<input type="checkbox"/>	0,839s	0,000s	Function	
<input checked="" type="checkbox"/> __srt_common_main	<input type="checkbox"/>	0,839s	0,000s	Function	
<input checked="" type="checkbox"/> __srt_common_main_seh	<input type="checkbox"/>	0,839s	0,000s	Function	
<input checked="" type="checkbox"/> invoke_main	<input type="checkbox"/>	0,839s	0,000s	Function	
<input checked="" type="checkbox"/> main	<input type="checkbox"/>	0,839s	0,000s	Function	
<input checked="" type="checkbox"/> [loop in mmult_task at mmult.cpp:131]	<input type="checkbox"/>	0,003s	0,000s	Scalar	
<input checked="" type="checkbox"/> mmult_task	<input type="checkbox"/>	0,839s	0,000s	Function	
<input checked="" type="checkbox"/> operator new[]	<input type="checkbox"/>	0,003s	0,000s	Function	
<input checked="" type="checkbox"/> operator new	<input type="checkbox"/>	0,003s	0,000s	Function	
<input checked="" type="checkbox"/> multiply_simple<float>	<input type="checkbox"/>	0,820s	97,8%	Function	

Bottom-up представление списка хотспотов приложения в виде таблицы с performance метриками.

Задание: Определить главный хотспот. Затем выделить ограничивающие его крыши, зафиксировать в отчете.

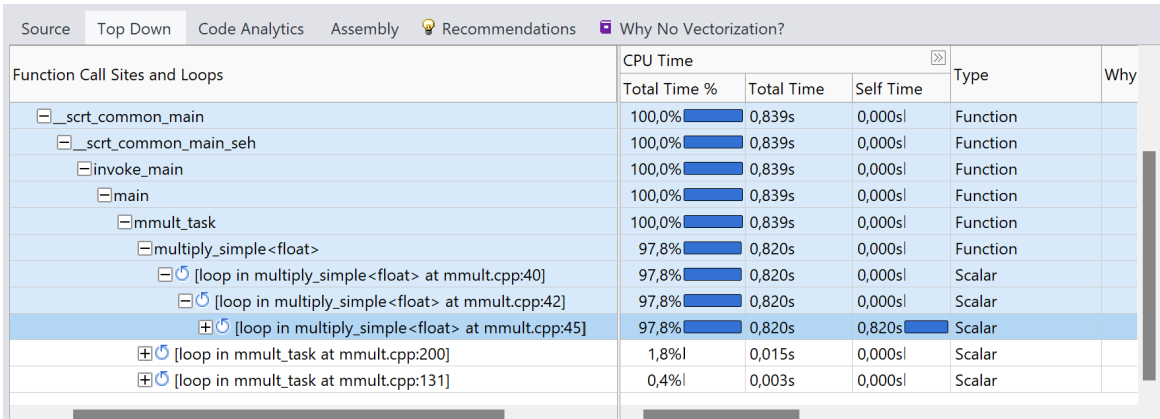
Ниже вкладки:

Source:



Выделено непосредственно место в исходном коде, соответствующее выбранному хотспоту.

Top Down



Дерево вызовов, для выбранного хотспота. Соответственно, для одного хотспота может быть выделено сразу несколько стеков вызовов. И тут уже можно посмотреть долю каждого в суммарном времени для хотспота.

**Code Analytics** - суммарная информация по хостпоту + по миксу инструкций для него (соотношение разных типов инструкций - memory vs compute)

Recommendations



Source
Top Down
Code Analytics
Assembly
Recommendations
Why No Vectorization?

### ! Potential underutilization of FMA instructions

Your current hardware supports the AVX2 instruction set architecture (ISA), which enables the use of fused multiply-add (FMA) instructions. Improve performance by utilizing FMA instructions.

#### Target the higher ISA

Although static analysis presumes the loop may benefit from FMA instructions available with the AVX2 or higher ISA, no FMA instructions executed for this loop. To fix: Use the following compiler options:

- `xCORE-AVX2` to compile for machines with and without AVX2 support
- `axCORE-AVX2` to compile for machines with AVX2 support only
- `xCOMMON-AVX512` to compile for machines with AVX-512 support only
- `axCOMMON-AVX512` to compile for machines with and without AVX-512 support

Note: the compiler options may vary depending on the CPU microarchitecture.

### ! Roofline conclusions

Conclusions, with optimization recommendations, are sorted by relevance.

#### This loop is mostly memory bound but may also be compute bound

The performance of the loop is bounded by the bandwidth of the shared cache and DRAM. To improve performance: Improve caching efficiency. The loop is also scalar. To fix: Vectorize the loop.

#### Collect Roofline for all memory levels

Run the [Roofline](#) for all memory levels to get a detailed analysis of memory-bound loops/functions.

Memory-Level Roofline evaluates the traffic between each memory subsystem based on cache simulation data.

### Potential underutilization of FMA instructions

Target the higher ISA

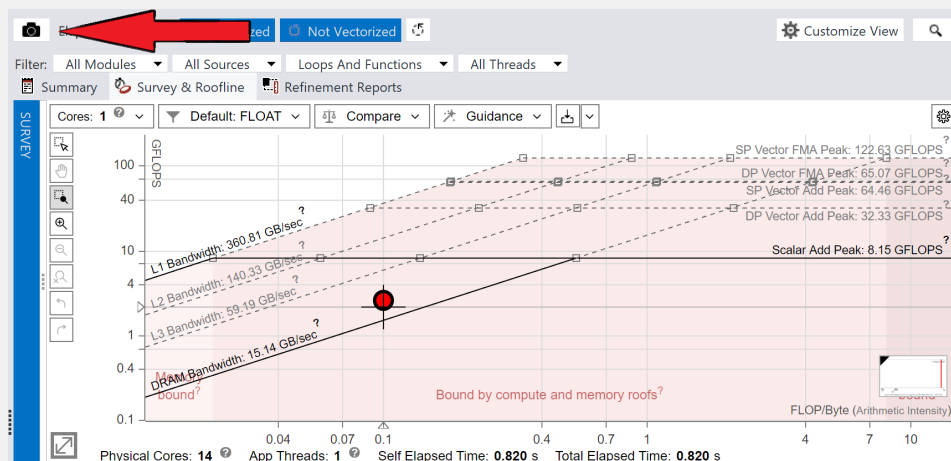
### Roofline conclusions

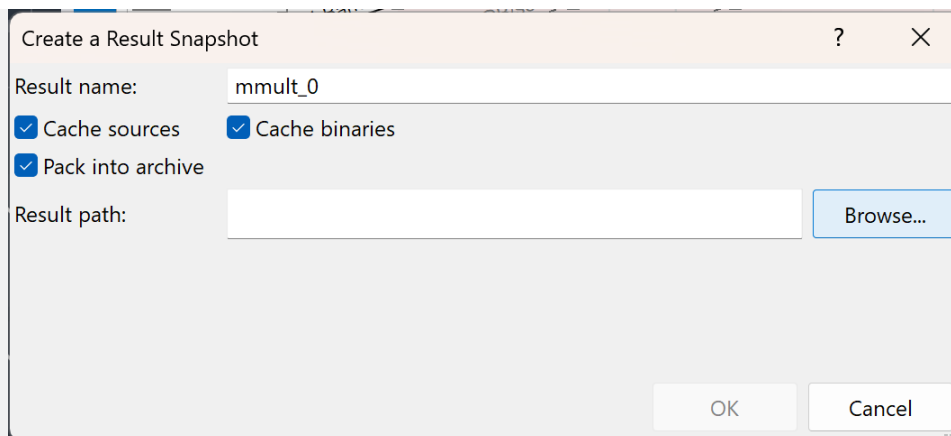
This loop is mostly memory bound but may also be compute bound

Collect Roofline for all memory levels

То, на что стоит обратить большее внимание. Довольно важной частью данного инструмента является система предоставления рекомендаций по анализируемому коду, советы, на что стоит обратить внимание, и даже, возможно, как можно исправить те или иные проблемы с производительностью приложения.

Далее перед любыми модификациями приложения необходимо делать снимки:





Снимки профиля, которые далее можно будет сравнивать между собой.

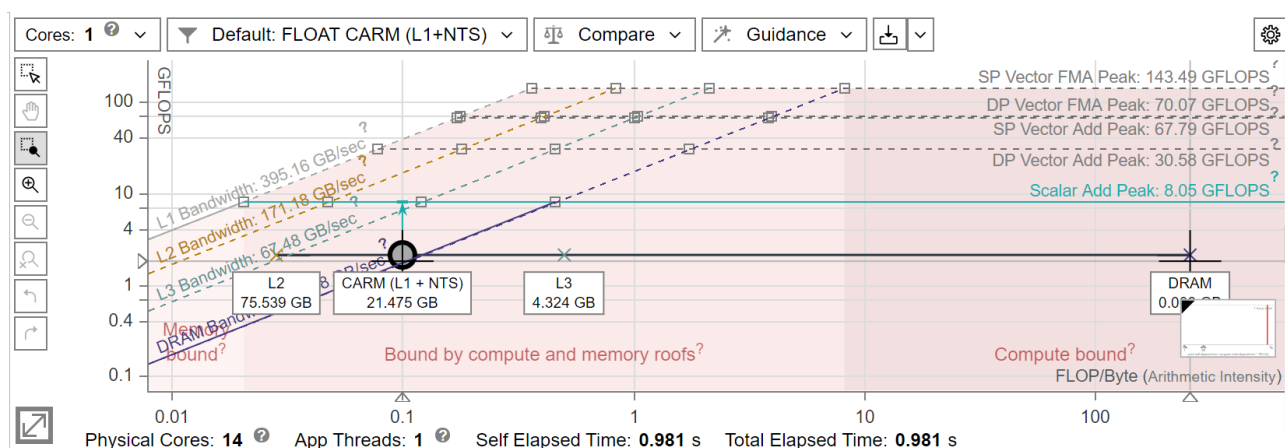
Итак, в данном примере нам предлагается две рекомендации:

- Использовать более широкий набор инструкций + FMA
- Собрать Roofline для всех уровней памяти.

Можете попытаться, в соответствии с рекомендациями, поправить в CMakeLists.txt для примера 1\_mmult используемое векторное расширение, результат зафиксировать.

В рамках этого разбора же попробуем собрать Memory Level Roofline.

Для этого необходимо в панели **Analysis Workflow**, в **Characterization** блоке выбрать пункт **Cache Simulation** и пересобрать отчет.



И теперь по двойному клику на точку на графике руфлайна раскроется раскладка для данного хотспота по утилизации различных уровней подсистемы памяти.

### ! Possible inefficient memory access patterns present

Inefficient memory access patterns may result in significant vector code execution slowdown or block automatic vectorization by the compiler. Improve performance by investigating.

#### Confirm inefficient memory access patterns

There is no confirmation inefficient memory access patterns are present. To fix: Run a [Memory Access Patterns analysis](#).

А в рекомендациях появится новая опция-предложение: запустить MAP анализ, который позволит определить паттерны доступа к памяти в рамках хотспота.

Для этого, во-первых, необходимо выделить интересующие нас хотспоты:

Function Call Sites and Loops	Performance Issues	CPU Time		Type	Why No Vectoriz
		Total Time	Self Time		
[loop in multiply_simple<float> at mmult.cpp:45]	2 Misaligned lo...	0,105s	0,105s	Inside vectorized	
mainCRTStartup		0,105s	0,000s	Function	
__schr_common_main		0,105s	0,000s	Function	
__schr_common_main_seh		0,105s	0,000s	Function	
invoke_main		0,105s	0,000s	Function	
main		0,105s	0,000s	Function	
mmult_task		0,105s	0,000s	Function	
multiply_simple<float>		0,105s	0,000s	Function	
[loop in multiply_simple<float> at mmult.cpp:40]		0,105s	0,000s	Vectorized (Body)	
[loop in multiply_simple<float> at mmult.cpp:40]		0,105s	0,000s	Scalar	

И выбрать в панели **Analysis Workflow** новый анализ - **Memory Access Patterns** и запустить его.

Третья вкладка **Refinement Reports**:

Тут можно найти отчет по проведенному MAP анализу:

Refinement Reports				
Site Location	Loop-Carried Dependencies	Strides Distribution	Access Pattern	Footprint Estimate
				Max. Per-Instruction Addr. R
[loop in multiply_simple<float> at mmult.cpp:40]	No Information Available	50% / 0% / 50%	Mixed Strides	448MB
[loop in multiply_simple<float> at mmult.cpp:42]	No Information Available	100% / 0% / 0%	All Unit Strides	447MB
[loop in multiply_simple<float> at mmult.cpp:45]	No Information Available	33% / 33% / 33%	Mixed Strides	447MB
<pre>43 { 44   #pragma vector 45   for (int k = 0; k &lt; arrSize; k++) 46   { 47       product[i][j] += aMatrix[i][k] * bMatrix[k][j]; 48   } 49 }</pre>				

Где, собственно, и написано, что в нашей программе, в хотспоте, обнаружен неэффективный паттерн доступа к памяти. А ниже приведены рокомендации, как от этого можно избавиться:

Memory Access Patterns Report   Dependencies Report   Recommendations

All Advisor-detectable issues: [C++](#) | [Fortran](#)

**Inefficient memory access patterns present**

There is a high of percentage memory instructions with irregular (variable or random) stride accesses. Improve performance by investigating and handling accordingly.

**Reorder loops**

This loop has less efficient memory access patterns than a nearby outer loop. To fix: Reorder the loops if possible.

Example (original code)

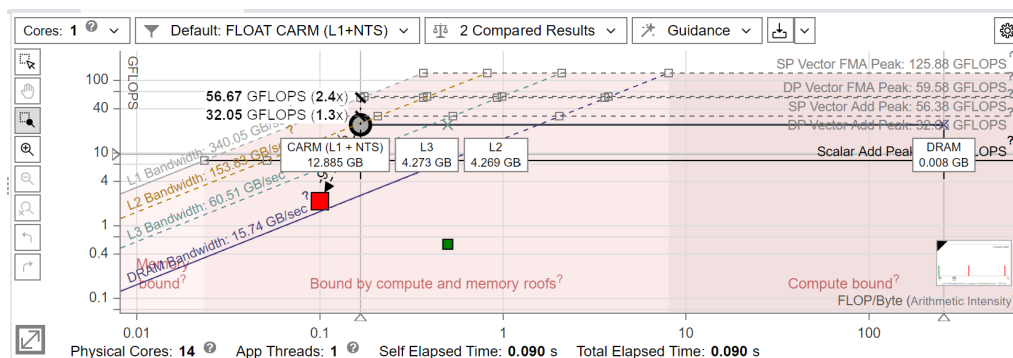
```
...
for (int j = 0; j < N; j++)
    for (int k = 0; k < N; k++)
        c[i][j] = c[i][j] + a[i][k] * b[k][j];
...
```

Example (revised code)

```
...
for (int k = 0; k < N; k++)
    for (int j = 0; j < N; j++)
        c[i][j] = c[i][j] + a[i][k] * b[k][j];
...
```

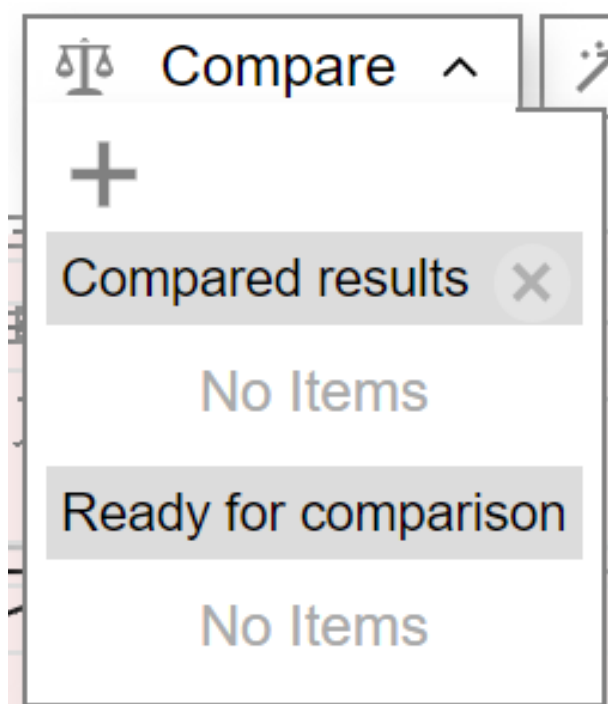
Сохраняем снэпшот и применяем исправления в коде. Снэпшот приложить к отчету. Можно собрать все в архив и сохранить для дальнейшей отправки преподавателю.

Memory Level Roofline для новой версии приложения:

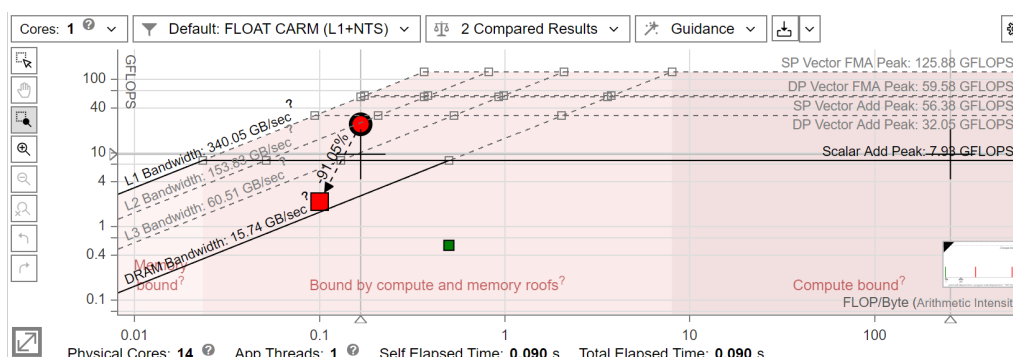


Доп. вопрос: о чем на рисунке выше говорит подобное расположение точек для разных уровней памяти на графике риффлайна.

А чтобы можно было сравнить две версии отчета соотв. разным версиям приложения:



Добавляем сохраненный ранее снимок и получаем следующую картинку:



Наглядно видно прирост производительности.

Для самостоятельного разбора остается два анализа: [Dependencies](#) и [Suitability](#)

Они помогут понять, что вообще возможно распараллелить, и насколько это будет эффективно масштабироваться.