

Recommendations:

Increase execution time:

Application execution time is too short. Metrics data may be unreliable. Consider reducing the sampling interval or increasing your application execution time.

Hotspots: Start with Hotspots analysis to understand the efficiency of your algorithm.

Use Hotspots analysis to identify the most time consuming functions. Drill down to see the time spent on every line of code.

Microarchitecture Exploration: There is low microarchitecture usage (15.0%) of available hardware resources. of Pipeline Slots

Run Microarchitecture Exploration analysis to analyze CPU microarchitecture bottlenecks that can affect application performance.

Memory Access: The Memory Bound metric is high (26.1%). A significant fraction of execution pipeline slots could be stalled due to demand memory load and stores. of Pipeline Slots

Use Memory Access analysis to measure metrics that can identify memory access issues.

Threading: There is poor utilization of logical CPU cores (69.3%) in your application.

Use Threading to explore more opportunities to increase parallelism in your application.

Elapsed Time: 0.048s

Application execution time is too short. Metrics data may be unreliable. Consider reducing the sampling interval or increasing your application execution time.

CPU:

IPC: 0.768

The IPC may be too low. This could be caused by issues such as memory stalls, instruction starvation, branch misprediction or long latency instructions. Explore the other hardware-related metrics to identify what is causing low IPC.

SP GFLOPS: 0.002

DP GFLOPS: 0.000

x87 GFLOPS: 0.001

Average CPU Frequency: 971.516 MHz

GPU:

Time: 37.5% (0.018s) of Elapsed time

GPU utilization is low. Consider offloading more work to the GPU to increase overall application performance.

IPC Rate: 1.286

Effective Logical Core Utilization: 69.3% (5.541 out of 8)

The metric value is low, which may signal a poor utilization of logical CPU cores while the utilization of physical cores is acceptable. Consider using logical cores, which in some cases can improve processor throughput and overall performance of multi-threaded applications.

Effective Physical Core Utilization: 100.0% (4.000 out of 4)

Microarchitecture Usage: 15.0% of Pipeline Slots

You code efficiency on this platform is too low.

Possible cause: memory stalls, instruction starvation, branch misprediction or long latency instructions.

Next steps: Run Microarchitecture Exploration analysis to identify the cause of the low microarchitecture usage efficiency.

Retiring: 15.0% of Pipeline Slots

Front-End Bound: 24.3% of Pipeline Slots

Issue: A significant portion of Pipeline Slots is remaining empty due to issues in the Front-End.

Tips: Make sure the code working size is not too large, the code layout does not require too many memory accesses

per cycle to get enough instructions for filling four pipeline slots, or check for microcode assists.

Back-End Bound:

54.9% of Pipeline Slots

A significant portion of pipeline slots are remaining empty. When operations take too long in the back-end, they introduce bubbles in the pipeline that ultimately cause fewer pipeline slots containing useful work to be retired per cycle than the machine is capable to support. This opportunity cost results in slower execution. Long-latency operations like divides and memory operations can cause this, as can too many operations being directed to a single execution port (for example, more multiply operations arriving in the back-end per cycle than the execution unit can support).

Memory Bound:

26.1% of Pipeline Slots

The metric value is high. This can indicate that the significant fraction of execution pipeline slots could be stalled due to demand memory load and stores. Use Memory Access analysis to have the metric breakdown by memory hierarchy, memory bandwidth information, correlation by memory objects.

Core Bound:

28.8% of Pipeline Slots

This metric represents how much Core non-memory issues were of a bottleneck. Shortage in hardware compute resources, or dependencies software's instructions are both categorized under Core Bound. Hence it may indicate the machine ran out of an OOO resources, certain execution units are overloaded or dependencies in program's data- or instruction- flow are limiting the performance (e.g. FP-chained long-latency arithmetic operations).

Bad Speculation:

5.8% of Pipeline Slots

Memory Bound:

26.1% of Pipeline Slots

The metric value is high. This can indicate that the significant fraction of execution pipeline slots could be stalled due to demand memory load and stores. Use Memory Access analysis to have the metric breakdown by memory

hierarchy, memory bandwidth information, correlation by memory objects.

L1 Bound: 28.4% of Clockticks

This metric shows how often machine was stalled without missing the L1 data cache. The L1 cache typically has the shortest latency. However, in certain cases like loads blocked on older stores, a load might suffer a high latency even though it is being satisfied by the L1.

L2 Bound: 0.3% of Clockticks

L3 Bound: 2.5% of Clockticks

DRAM Bound: 5.6% of Clockticks

DRAM Bandwidth Bound: 0.0% of Elapsed Time

Store Bound: 0.7% of Clockticks

Vectorization: 0.0% of Packed FP Operations

A significant fraction of floating point arithmetic instructions are scalar. Use Intel Advisor to see possible reasons why the code was not vectorized.

Instruction Mix:

SP FLOPs: 0.1% of uOps

Packed: 0.0% from SP FP

128-bit: 0.0% from SP FP

256-bit: 0.0% from SP FP

Scalar: 100.0% from SP FP

This code has floating point operations and is not vectorized. Consider using Intel Advisor to vectorize the loops.

DP FLOPs: 0.0% of uOps

Packed: 4.2% from DP FP

128-bit: 4.2% from DP FP

256-bit: 0.0% from DP FP

Scalar: 95.8% from DP FP

A significant fraction of floating point arithmetic instructions are scalar. Use Intel Advisor to see possible reasons why the code was not vectorized.

x87 FLOPs: 0.0% of uOps

Non-FP: 99.9% of uOps

FP Arith/Mem Rd Instr. Ratio: 0.002

The metric value is low. This can be a result of unaligned access to data for vector operations. Use Intel Advisor to find possible data access inefficiencies for vector operations.

FP Arith/Mem Wr Instr. Ratio: 0.005

The metric value is low. This can be a result of unaligned access to data for vector operations. Use Intel Advisor to find possible data access inefficiencies for vector operations.

GPU Active Time: 37.5%

GPU utilization is low. Consider offloading more work to the GPU to increase overall application performance.

GPU Utilization when Busy: 17.8%

The percentage of time when the EUs were stalled or idle is high, which has a negative impact on compute-bound applications.

IPC Rate: 1.286

EU State: 17.8%

Active: 17.8%

Stalled: 33.8%

A significant portion of GPU time is lost due to stalls. For compute-bound code, this could indicate that performance is limited by memory or sampler accesses.

Idle: 48.4%

A significant portion of GPU time is spent idle. This is usually caused by imbalance or thread scheduling problems.

Occupancy: 31.4% of peak value

Low value of the occupancy metric may be caused by inefficient work scheduling. Make sure work items are neither too small nor too large.

Collection and Platform Info:

Application Command Line: ./codecs/HHI-VVC/decoder/vvdecapp "-b" ".bin/HHI-VVC/randomaccess_fast.cfg/CLASS_C/RaceHorses_416x240_30_QP_27_HHI-VVC.bin"

Operating System: 5.4.0-72-generic DISTRIB_ID=Ubuntu
DISTRIB_RELEASE=18.04 DISTRIB_CODENAME=bionic
DISTRIB_DESCRIPTION="Ubuntu 18.04.5 LTS"

Computer Name: eimon

Result Size: 3.7 MB

Collection start time: 22:29:28 18/04/2021 UTC

Collection stop time: 22:29:28 18/04/2021 UTC

Collector Type: Event-based sampling driver, Event-based counting driver

CPU:

Name: Intel(R) Processor code named Kabylake ULX

Frequency: 1.992 GHz

Logical CPU Count: 8

Max DRAM Single-Package Bandwidth: 11.000 GB/s

Cache Allocation Technology:

Level 2 capability: not detected

Level 3 capability: not detected

GPU:

Name:	Display controller: Intel Corporation Device 22807
Vendor:	Intel Corporation
EU Count:	24
Max EU Thread Count:	7
Max Core Frequency:	1.150 GHz

Intel® oneAPI VTune™ Profiler 2021.1.1 Gold**Recommendations:****Increase execution time:**

Application execution time is too short. Metrics data may be unreliable. Consider reducing the sampling interval or increasing your application execution time.

Hotspots: Start with Hotspots analysis to understand the efficiency of your algorithm.

Use Hotspots analysis to identify the most time consuming functions. Drill down to see the time spent on every line of code.

Microarchitecture Exploration: There is low microarchitecture usage (1.7%) of available hardware resources. of Pipeline Slots

Run Microarchitecture Exploration analysis to analyze CPU microarchitecture bottlenecks that can affect application performance.

Memory Access: The Memory Bound metric is high (48.9%). A significant fraction of execution pipeline slots could be stalled due to demand memory load and stores. of Pipeline Slots

Use Memory Access analysis to measure metrics that can identify memory access issues.

Elapsed Time: 0.028s

Application execution time is too short. Metrics data may be unreliable. Consider reducing the sampling interval or increasing your application execution time.

CPU:**IPC:** 0.056

The IPC may be too low. This could be caused by issues such as memory stalls, instruction starvation, branch misprediction or long latency instructions. Explore the other hardware-related metrics to identify what is causing low IPC.

DP GFLOPS: 0.000

x87 GFLOPS: 0.001

Average CPU Frequency: 1.988 GHz

GPU:

Time: 31.2% (0.009s) of Elapsed time

GPU utilization is low. Consider offloading more work to the GPU to increase overall application performance.

IPC Rate: 1.319

Effective Logical Core Utilization: 132.9% (10.634 out of 8)

Effective Physical Core Utilization: 100.0% (4.000 out of 4)

Microarchitecture Usage: 1.7% of Pipeline Slots

You code efficiency on this platform is too low.

Possible cause: memory stalls, instruction starvation, branch misprediction or long latency instructions.

Next steps: Run Microarchitecture Exploration analysis to identify the cause of the low microarchitecture usage efficiency.

Retiring: 1.7% of Pipeline Slots

Front-End Bound: 8.9% of Pipeline Slots

Back-End Bound: 88.1% of Pipeline Slots

A significant portion of pipeline slots are remaining empty. When operations take too long in the back-end, they introduce bubbles in the pipeline that ultimately cause fewer pipeline slots containing useful work to be retired per cycle than the machine is capable to support. This opportunity cost results in slower execution. Long-latency

operations like divides and memory operations can cause this, as can too many operations being directed to a single execution port (for example, more multiply operations arriving in the back-end per cycle than the execution unit can support).

Memory Bound:

48.9% of Pipeline Slots

The metric value is high. This can indicate that the significant fraction of execution pipeline slots could be stalled due to demand memory load and stores. Use Memory Access analysis to have the metric breakdown by memory hierarchy, memory bandwidth information, correlation by memory objects.

Core Bound:

39.2% of Pipeline Slots

This metric represents how much Core non-memory issues were of a bottleneck. Shortage in hardware compute resources, or dependencies software's instructions are both categorized under Core Bound. Hence it may indicate the machine ran out of an OOO resources, certain execution units are overloaded or dependencies in program's data- or instruction- flow are limiting the performance (e.g. FP-chained long-latency arithmetic operations).

Bad Speculation:

1.3% of Pipeline Slots

Memory Bound:

48.9% of Pipeline Slots

The metric value is high. This can indicate that the significant fraction of execution pipeline slots could be stalled due to demand memory load and stores. Use Memory Access analysis to have the metric breakdown by memory hierarchy, memory bandwidth information, correlation by memory objects.

L1 Bound:

36.1% of Clockticks

This metric shows how often machine was stalled without missing the L1 data cache. The L1 cache typically has the shortest latency. However, in certain cases like loads blocked on older stores, a load might suffer a high latency even though it is being satisfied by the L1.

L2 Bound: 0.4% of Clockticks
L3 Bound: 2.6% of Clockticks
DRAM Bound: 5.9% of Clockticks
DRAM Bandwidth Bound: 70.5% of Elapsed Time

The system spent much time heavily utilizing DRAM bandwidth. Improve data accesses to reduce cacheline transfers from/to memory using these possible techniques: 1) consume all bytes of each cacheline before it is evicted (for example, reorder structure elements and split non-hot ones); 2) merge compute-limited and bandwidth-limited loops; 3) use NUMA optimizations on a multi-socket system. Note: software prefetches do not help a bandwidth-limited application. Run Memory Access analysis to identify data structures to be allocated in High Bandwidth Memory (HBM), if available.

Store Bound: 0.2% of Clockticks

Vectorization: 0.0% of Packed FP Operations

A significant fraction of floating point arithmetic instructions are scalar. Use Intel Advisor to see possible reasons why the code was not vectorized.

Instruction Mix:

SP FLOPs:	0.0% of uOps
Packed:	0.0% from SP FP
128-bit:	0.0% from SP FP
256-bit:	0.0% from SP FP
Scalar:	0.0% from SP FP
DP FLOPs:	0.0% of uOps
Packed:	0.0% from DP FP
128-bit:	0.0% from DP FP
256-bit:	0.0% from DP FP
Scalar:	100.0% from DP FP

This code has floating point operations and is not vectorized. Consider using Intel Advisor to vectorize the loops.

x87 FLOPs: 0.1% of uOps

Non-FP: 99.9% of uOps

FP Arith/Mem Rd Instr. Ratio: 0.004

The metric value is low. This can be a result of unaligned access to data for vector operations. Use Intel Advisor to find possible data access inefficiencies for vector operations.

FP Arith/Mem Wr Instr. Ratio: 0.006

The metric value is low. This can be a result of unaligned access to data for vector operations. Use Intel Advisor to find possible data access inefficiencies for vector operations.

GPU Active Time: 31.2%

GPU utilization is low. Consider offloading more work to the GPU to increase overall application performance.

GPU Utilization when Busy: 22.4%

The percentage of time when the EUs were stalled or idle is high, which has a negative impact on compute-bound applications.

IPC Rate: 1.319
EU State: 22.4%
Active: 22.4%
Stalled: 35.3%

A significant portion of GPU time is lost due to stalls. For compute-bound code, this could indicate that performance is limited by memory or sampler accesses.

Idle: 42.2%

A significant portion of GPU time is spent idle. This is usually caused by imbalance or thread scheduling problems.

Occupancy: 36.5% of peak value

Low value of the occupancy metric may be caused by inefficient work scheduling. Make sure work items are neither too small nor too large.

Collection and Platform Info:

Application Command Line: ./codecs/HHI-VVC/decoder/vvdecapp "-b" ".bin/HHI-VVC/randomaccess_fast.cfg/CLASS_C/RaceHorses_416x240_30_QP_27_HHI-VVC.bin"

Operating System: 5.4.0-72-generic DISTRIB_ID=Ubuntu DISTRIB_RELEASE=18.04 DISTRIB_CODENAME=bionic DISTRIB_DESCRIPTION="Ubuntu 18.04.5 LTS"

Computer Name: eimon

Result Size: 3.7 MB

Collection start time: 07:49:15 19/04/2021 UTC

Collection stop time: 07:49:15 19/04/2021 UTC

Collector Type: Event-based sampling driver,Event-based counting driver

CPU:

Name: Intel(R) Processor code named Kabylake ULX

Frequency: 1.992 GHz

Logical CPU Count: 8

Max DRAM Single-Package Bandwidth: 10.000 GB/s

Cache Allocation Technology:

Level 2 capability: not detected

Level 3 capability: not detected

GPU:

Name: Display controller: Intel Corporation Device 22807

Vendor:	Intel Corporation
EU Count:	24
Max EU Thread Count:	7
Max Core Frequency:	1.150 GHz