

NPU validation & debug

May 2024

CCG CPE CCE

Phoebe Hsu

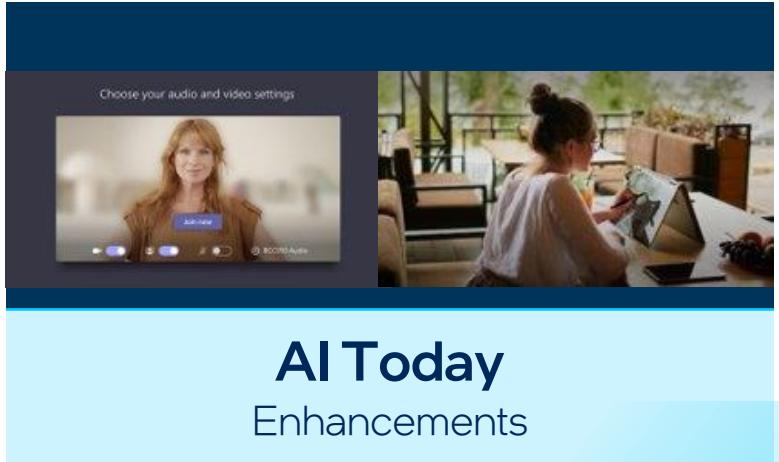


Outline

- Client AI and Roadmap
- Validation
 - Workshop validation
 - NPU benchmarking
 - Procyon (3rd party app)
 - PTAT
 - MSFT WSE Tool
- Logs/Debug
 - Intel STC v2.3: NPU/MEP Camera trace/UMD/KMD, Xperf
 - GPUView
 - Windbg-Live kernel dump/BSOD memory dump
 - Traceview
 - WPR
 - PTAT

Client AI and Roadmap

Transforming the PC Experience



Elevated video collaboration & streaming
Enhanced Audio effects
Creator and Gaming effects



AI Tomorrow
Everything

AI Assistants know your daily context
More creative, productive, & collaborative
Across everything you do

Cloud

Massive scalable compute
High Latency
Privacy Concerns
Expensive

Client

Massive distributed scale
Low Latency
Improved Privacy
Lower Cost (to ISV)

8Q Client AI Roadmap

| | CPU | iGPU | NPU | |
|---------------|----------------------------|--|--|--|
| | Fast Response, Low latency | AI integrated with 3D/render/ media pipelines; high batch size | Sustained Workloads, dedicated offload | |
| | Meteor Lake | Lunar Lake (MX 17W) | Arrow Lake | Panther Lake |
| CPU | 6 (45W) U-series | 5 H-series | 9 (45W) H-series | 14 S-series |
| iGPU | 8 | 19 | 59 | 72 U-series |
| NPU | 11 | 45 | 11 | 9 H-series |
| Platform TOPS | 25 | 36 | 109 | Up to 92 Up to 36 98 182 |

8Q Client AI Roadmap (Cont.)

Client AI Workloads are Diverse
No Single Compute Unit Meets All Key Needs



Bursty, Latency sensitive



Sustained, Battery life sensitive



Periodic, Throughput sensitive

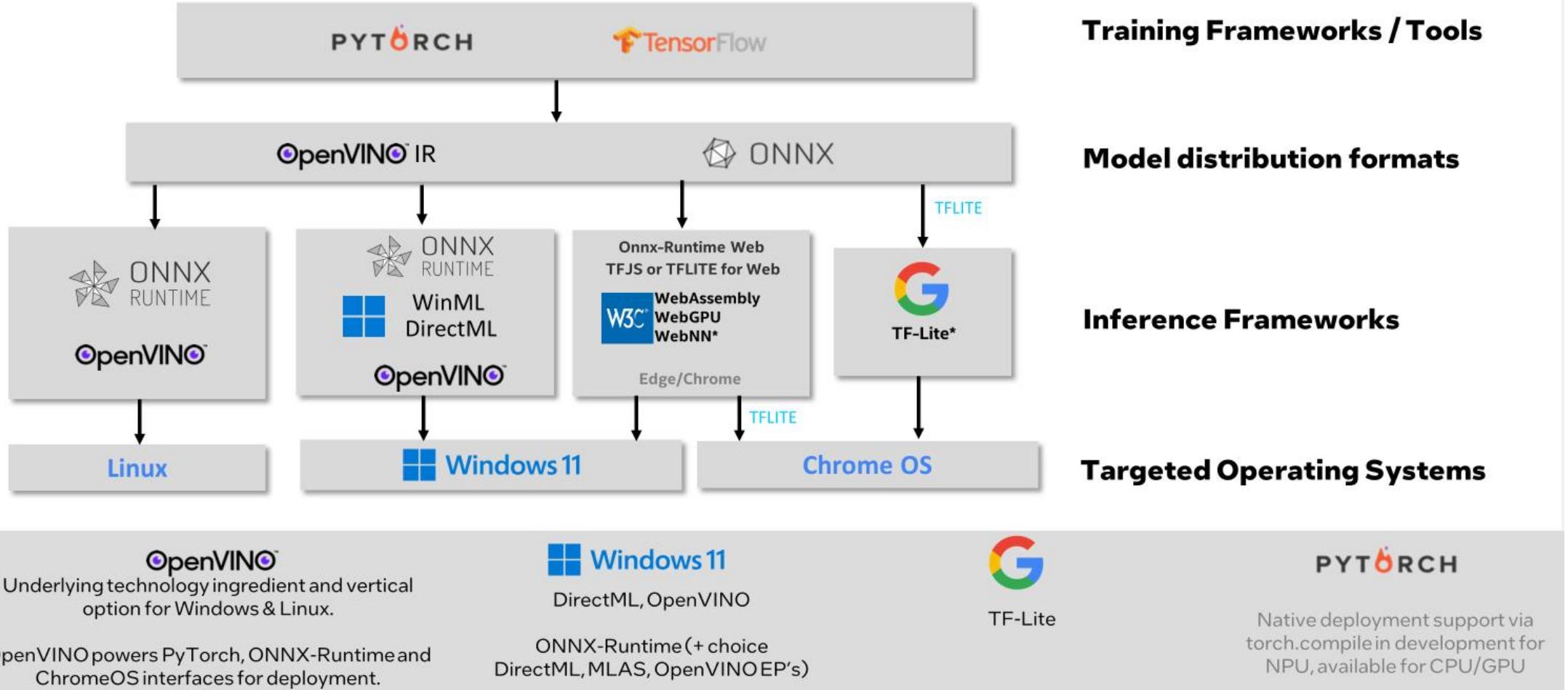
| HW | Value | RPL | MTL | ARL | LNL MX | PTL |
|------|---|---|---|--|---|--|
| CPU | SW Programmability; low latency, single inference tasks | AVX-256 VNNI H: 4-5 TOPS | AVX-256 VNNI H: ~3-6; U: ~2-3 TOPS | AVX-256 VNNI H: ~7-9.5; S: 14 TOPS | AVX-256 VNNI ~2-5 TOPS | AVX2+ TOPS - H: Up to 11; U: 5 |
| iGPU | AI integrated with 3D/render/media pipelines; high batch size | DP4A H/U: up to 9 TOPS S/HX: 3 TOPS | DP4a(U, H) H: up to 19 TOPS U: up to 8 TOPS | DP4a(U, S, HX) ~9 TOPS ARL H w/X ^e Matrix Extensions (XMX) Up to 72 TOPS | DP4a + X ^e Matix Extensions (XMX) Up to 59 TOPS | DP4a + XMX H: Up to ~123 TOPS U: up to 41 TOPS |
| iNPU | Dedicated AI Offload, Power efficiency for Battery Life | NA | | NPU 2.7 TOPS - H: 11 TOPS; U: 9.5-11; ARL S, HX: 13 | NPU 4.0 Up to 45 TOPS | NPU 5.0 Up to 48 TOPS |

TOPS will vary slightly based on power & frequency of each sku

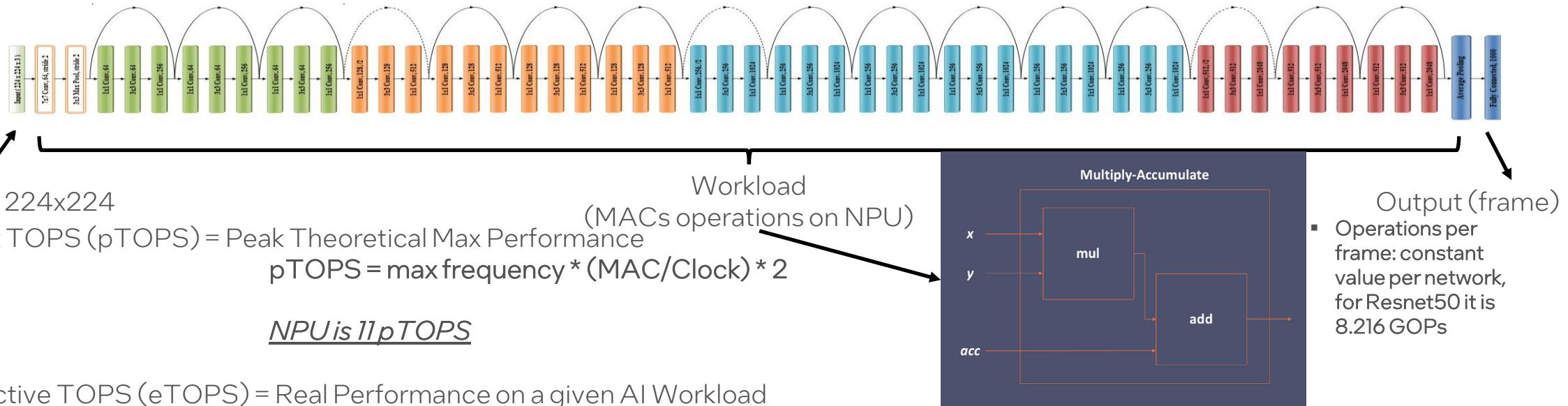
The Right Frameworks for Innovation and Scale:



Client AI Inference SW For Deployment 2023+



Resnet50 Example Based on MTL



- Peak TOPS (pTOPS) = Peak Theoretical Max Performance
 $pTOPS = \text{max frequency} * (\text{MAC/Clock}) * 2$

NPU is 11 pTOPS

- Effective TOPS (eTOPS) = Real Performance on a given AI Workload (the efficiency of pTOPS)

$$eTOPS = (\text{fps} * \text{each frame GOPs}) / 1000$$

$$\underline{\text{NPU } 8.2 \text{ eTOPS} = (1000 * 8.216) / 1000}$$

- We use ResNet50: a common network + a good mix of a memory and compute bound network. Is it Perfect? -> No, but it is better than pTOPS as eTOPS shows real workload measured across many HW configs
- AI Benchmark for Client: Not 1 standard Today (UL Procyon Redowa (POR)/MLPerf/GeekBenchML)

| | pTOPS | Resnet50 FPS | eTOPS | Efficiency |
|---------|-------|--------------|-------|------------|
| MTL NPU | 11 | 1000 | 8.2 | 75% |

Validation

Workshop validation

- Download [Intel NPU OEM Workshop Validation Software Kits](#)
 - \validation_mtl_arl_Inl\validation\test_scripts

```
benchmark.bat  
benchmark_latency.json  
benchmark_throughput.json  
classification.bat  
NPU_001.bat  
NPU_002.bat  
NPU_003.bat
```

- More MSFT MEP test case
 - [Meteor Lake VPU MEP Enablement and Validation Testcases in OEM Platform](#)
 - Hot-plug AC adaptor/type C monitor, MS/S4/Reboot stress test, screen on idle, teams call, launch/close camera app, on/off effects...etc.

Workshop validation

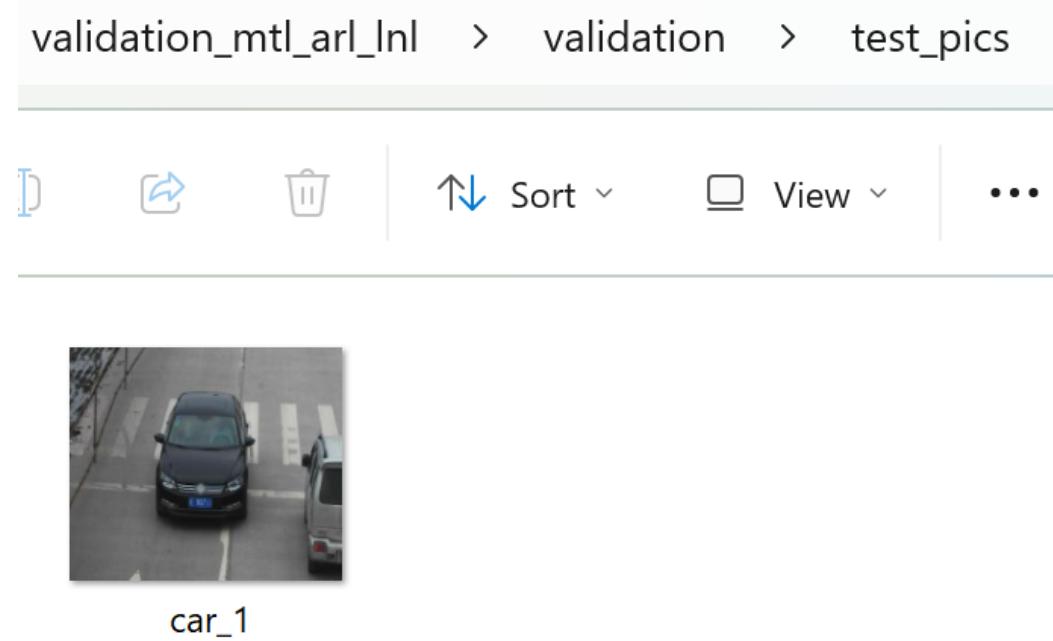
| | | | | |
|---------|-----|---|---|---|
| NPU_001 | NPU | Verify NPU with multiple inference models support using OpenVINO framework (ResNet/MobileNet) | Verify NPU with multiple inference models support using OpenVINO framework. The models will be run sequentially using the application. A single instance of Image Classification will be executed using ResNet followed by MobileNet. | a) Command runs successfully. No errors, BSODs, hangs, TDRs and crashes are observed. b) Open window showing tiles of classified image and the expectation is that the prediction of the input is correct. |
| NPU_002 | NPU | Verify NPU driver for D0 and D3 states | Verify NPU driver for D0 and D3 states. The OpenVINO model used is ResNet 50 and the application used to run the model on VPU IP is Image Classification. The transition from D3 to D0 and D0 to D3 is checked using SocWatch tool | a) Open SoCWatch without error, SoC Watch command : "socwatch.exe -f sys -o filename" b) User should be able to navigate to the mentioned path in device manager and note the current power state from it. If the NPU device is not working Power state should be D3. If the NPU device is working, the power state should be D0. c) Run inference and Stop log collection & verify VPU device should go to D3->D0->D3 |
| NPU_003 | NPU | Verify NPU functionality in DC mode real battery | Verify offloading workload to NPU with SUT connected to real battery | a) NPU utilization is observed to increase. NPU Compute utilization > 0% when inferencing is ongoing. b) Command runs successfully. No errors, BSODs, hangs, TDRs and crashes are observed. c) The first inference time, throughput and latency results are non zero value |
| NPU_004 | NPU | Verify the MEP camera effect on NPU | Verify the MEP installation and optin on NPU | a) Modify the ivd64extn.inf in NPU driver, default SUBSYS HWID LNL (8086_643E), MTL (8086_7D1D) b) Install the ivd64extn.inf and check the status in Device Manager c) Install MEP base package and camera component package, then check the status in Device Manager Software components, Windows Studio Effects should be available d) Optin the camera with CameraOptinUtils from MSFT, or add "FSMEnableMsEffects: REG_DWORD: 0x1" under Device Interface Node registry e) If the MEP and Camera Opt-in are successful, control the MEP camera effects in Windows Setting and test it with Windows Camera application |

Workshop validation

```
Image C:\openvino\test_scripts\..\test_pics\car_1.bmp

classid probability
-----
656    12.6250000
734    10.9765625
705    10.8203125
468    10.4453125
654    10.1640625
436    9.7421875
874    9.3281250
407    9.1171875
829    9.1562500
675    8.9140625

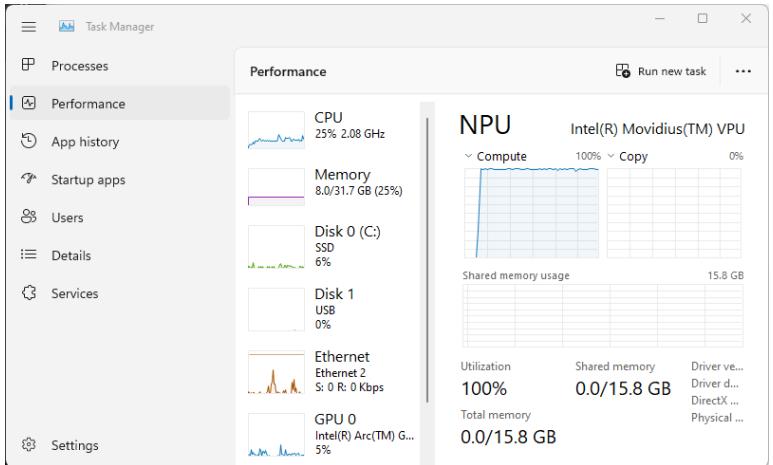
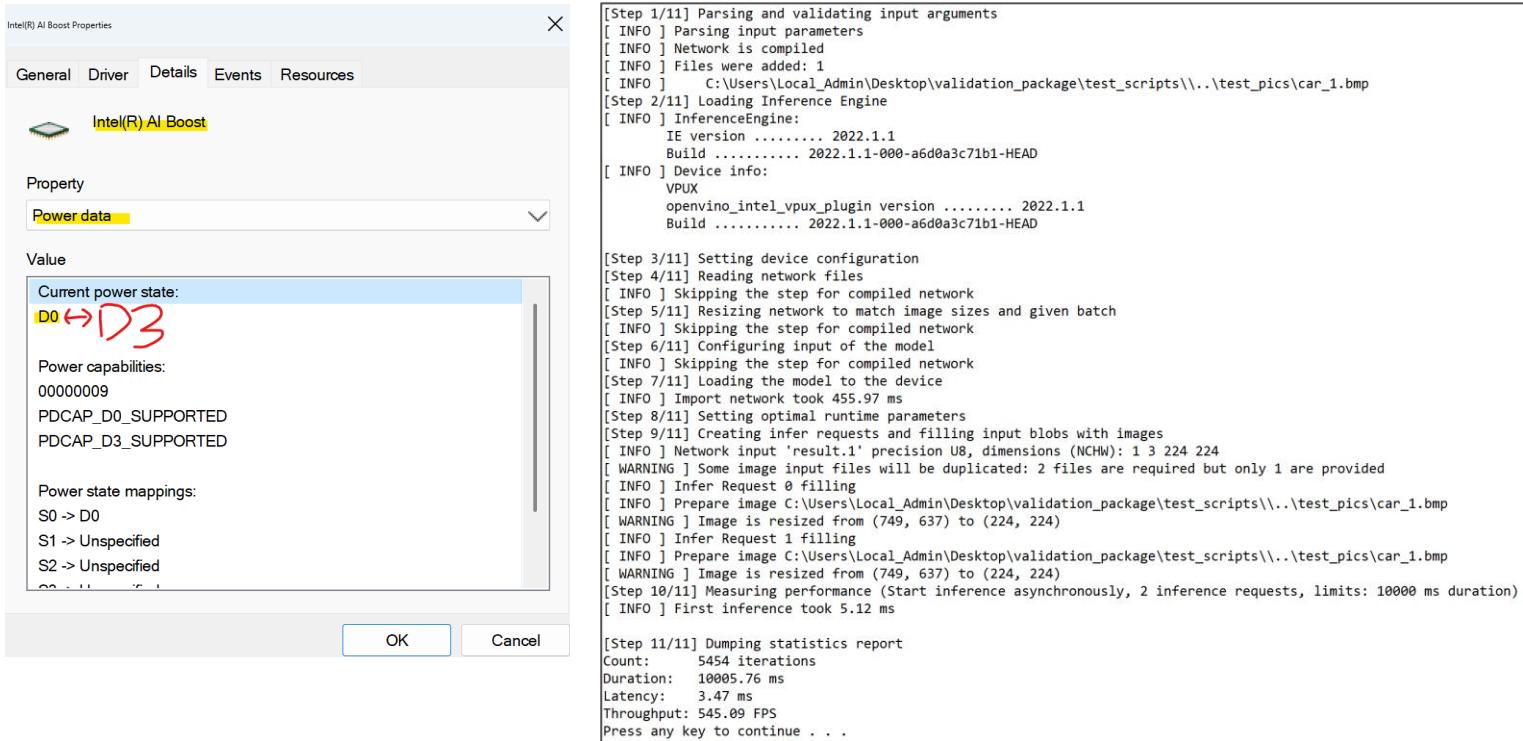
[ INFO ] Execution successful
```



656: 'minivan',
734: 'police van, police wagon, paddy wagon, patrol wagon, wagon, black Maria'
705: 'passenger car, coach, carriage'

1000 class ids to human readable labels
<https://gist.github.com/yrevar/942d3a0ac09ec9e5eb3a>

Workshop validation



NPU benchmarking

- We use ResNet50 to measure eTOPS
 - a common network + a good mix of a memory & compute bound network
- Effective TOPS (eTOPS) = real performance on a given AI workload (the efficiency of pTOPS)
 - $eTOPS = (\text{fps} * \text{each frame GOPs}) / 1000$
 - Operations per frame: constant value per network, for Resnet50 it's 8.216 GOPs
- See [MTL/LNL Platform NPU Benchmarking](#) for details

Run Resnet50 on MTL NPU

1. Download NPU driver from RDC and then install it and its extension driver w/ 4-part ID
2. Install one of Python 3.8 - 3.11 and following below steps to create virtual environment, ov_24.0 under C:\Users\Public\

```
$ cd C:\Users\Public\  
$ python -m venv ov_24.0  
$ ov_24.0\Scripts\activate  
$ python -m pip install --upgrade pip  
$ pip install openvino==2024.0
```

3. Download OpenVINO package, [w_openvino_toolkit_mtl_23ww19.zip](#), form WW22'23 BKC release to get Resnet50 model, and then unzip to C:\Users\Public, and you can see below 2 zip files

| | | | |
|--|---|------------------|------------|
| w_ov_dyn_win_t_230508_1951_ov_1f790df33c73a9a_vx_a9f38da07e3a135.zip | ✓ | 6/2/2023 3:27 PM | Compressed |
| w_ov_dyn_win_t_230508_1951_ov_1f790df33c73a9a_vx_a9f38da07e3a135_TOOLS.zip | ✓ | 6/2/2023 3:27 PM | Compressed |

4. Unzip "w_ov_dyn_win_t_230508_1951_ov_1f790df33c73a9a_vx_a9f38da07e3a135_TOOLS.zip" to C:\Users\Public
5. Create npu_config.json and add below line to it, and then save it at C:\Users\Public\

```
{"NPU": {"NPU_COMPILER_TYPE": "DRIVER", "NPU_COMPILATION_MODE_PARAMS": "enable-activation-sparsity=true" }}
```

6. Pre-check:

| Required Files | File Path |
|---|---|
| ResNet50 int8 with 50% sparsity (IR model) | <ul style="list-style-type: none">C:\Users\Public\tools\ir_models\IRv1\20230309_vpu-models-mtl-por-ir_v1\ov_2022.3.0-9752faf8eb\resnet-50-v1_5-sparse50\onnx\FP16-INT8\resnet-50-v1_5-sparse50.xmlC:\Users\Public\tools\ir_models\IRv1\20230309_vpu-models-mtl-por-ir_v1\ov_2022.3.0-9752faf8eb\resnet-50-v1_5-sparse50\onnx\FP16-INT8\resnet-50-v1_5-sparse50.bin |
| NPU configuration | C:\Users\Public\npu_config.json |

Run Resnet50 on MTL NPU (Cont.)

7. Open cmd terminal and then run below commands to setup environment and prepare model before run benchmark_app

```
$ C:\Users\Public\ov_24.0\Scripts\activate  
$ cd C:\Users\Public\  
$ copy C:\Users\Public\tools\ir_models\IRv1\20230309_vpu-models-mtl-por-ir_v11_ov_2022.3.0-9752faf8eb\resnet-50-v1_5-sparse50\onnx\FP16-  
INT8\resnet-50-v1_5-sparse50.*.
```

8. Run NPU in Throughput mode

```
$ benchmark_app -d NPU -m resnet-50-v1_5-sparse50.xml -t 180 -layout [NCHW] -ip f16 -op f16 -hint throughput -load_config npu_config.json
```

9. Run NPU in Latency mode

```
$ benchmark_app -d NPU -m resnet-50-v1_5-sparse50.xml -t 180 -layout [NCHW] -ip f16 -op f16 -api sync -load_config npu_config.json
```

Result for reference:

- Configuration:
 - MTL-H 28W
 - NPU 2196 driver
 - OpenVINO 2024.0
 - MEMORY size: 16GB
- Performance on MTL NPU:

| Execution Mode | NPU |
|----------------|------|
| Latency | 6.43 |
| Throughput | 8.47 |

Unit: eTOPs

One-click Script for Test

- Pre-check:

| Required Files | File Path |
|---|---|
| Python environment | C:\Users\Public\ov_24.0\ |
| ResNet50 int8 with 50% sparsity (IR model) | <ul style="list-style-type: none">• C:\Users\Public\resnet-50-v1_5-sparse50.xml• C:\Users\Public\resnet-50-v1_5-sparse50.bin |
| NPU configuration | C:\Users\Public\npu_config.json |

- Create resnet50.bat and add below lines into it, and then save it at C:\Users\Public\

```
call C:\Users\Public\ov_24.0\Scripts\activate
cd C:\Users\Public\
benchmark_app -d NPU -m resnet-50-v1_5-sparse50.xml -t 180 -ip f16 -op f16 -layout [NCHW] -api sync -load_config npu_config.json
benchmark_app -d NPU -m resnet-50-v1_5-sparse50.xml -t 180 -ip f16 -op f16 -layout [NCHW] -hint throughput -load_config npu_config.json

pause
```

Procyon

- Apply access

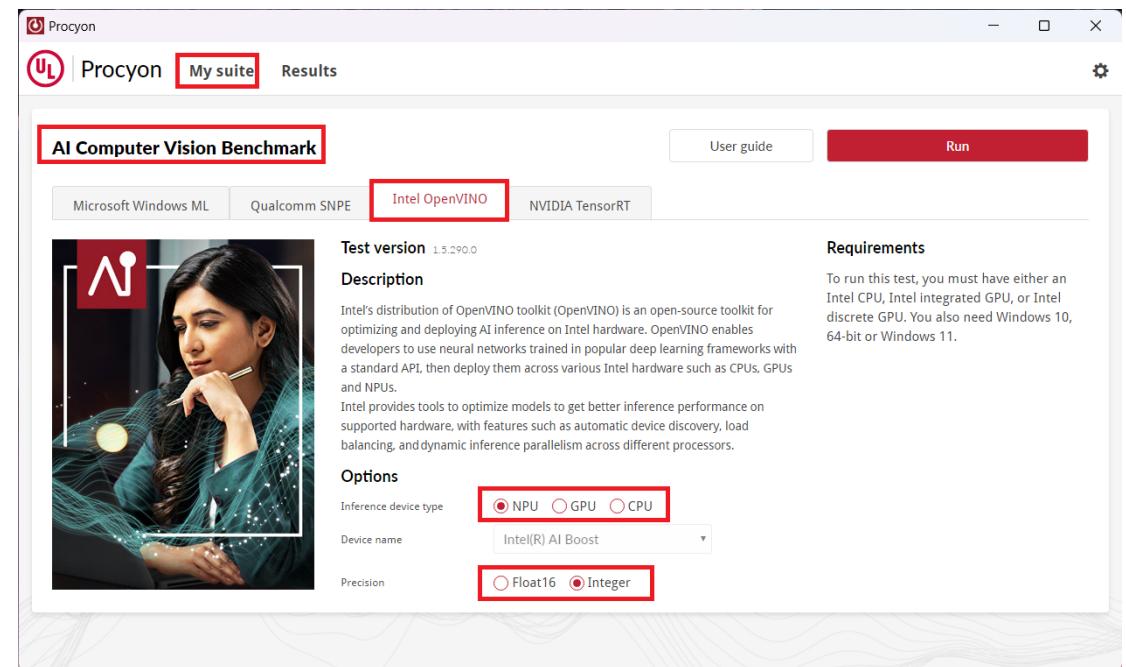
<https://wiki.ith.intel.com/display/ITSVPUCSE/How+to+get+Redowa+packages>

- Procyon package

<\\pwb-release.amr.corp.intel.com\release\Futuremark\Procyon> (Release & Pre-release)

- Procyon system configuration

- Procyon version: $\geq 2.6.896$
- NPU version: 31.0.100.1688
- Score: ≥ 480 (+- 5%)
- AC, Best Performance
- NPU, select OpenVino tab



Procyon

Example on MTL RVP

Procyon

Procyon My suite Results

Overall score **506**
4% higher than the previous result

CPU: Intel Core Ultra 7 1002H
GPU: Intel(R) Arc(TM) Graphics

Application versions
Intel OpenVINO version 2023.1.0

Settings used
Inference device used: Intel(R) NPU (3720VE)
Inference device type: NPU
Precision: integer



Detailed scores

MobileNet V3
Average inference time 0.71 ms Median inference time 0.69 ms Total inferences count 183 252

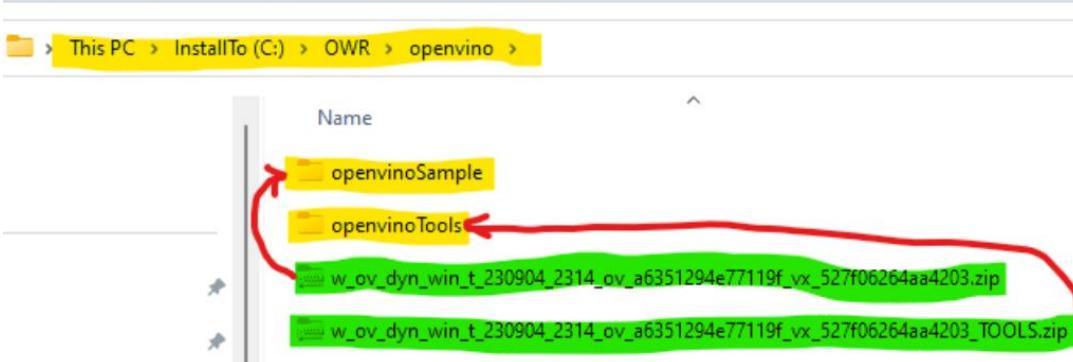
ResNet 50
Average inference time 1.68 ms Median inference time 1.66 ms Total inferences count 94 465

Inception V4

Back Save to cloud Show log Export as File Export as PDF Export as XML

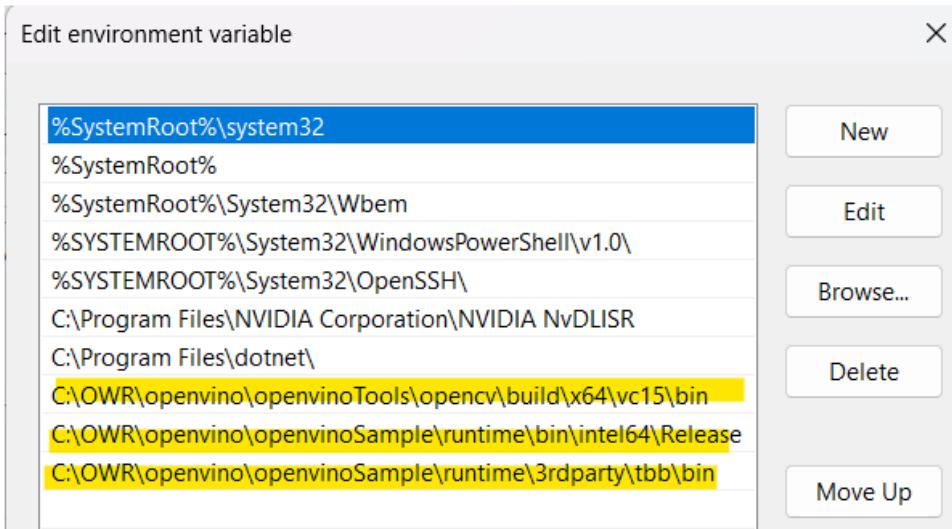
PTAT

- Download [Intel® Power And Thermal Analysis Tool](#)
- Download the latest toolkit “w_openvino_toolkit_YYWWNN”.zip (internal)
- Follow below steps to make the OpenVINO library files available for running the NPU workload
 - There are two w_vo_*.zip files of which one has the sample files, and another has tool files.
 - Unzip the sample zip file and tool zip file to a directory and rename it to “openvinoSample” and “openvinoTools”



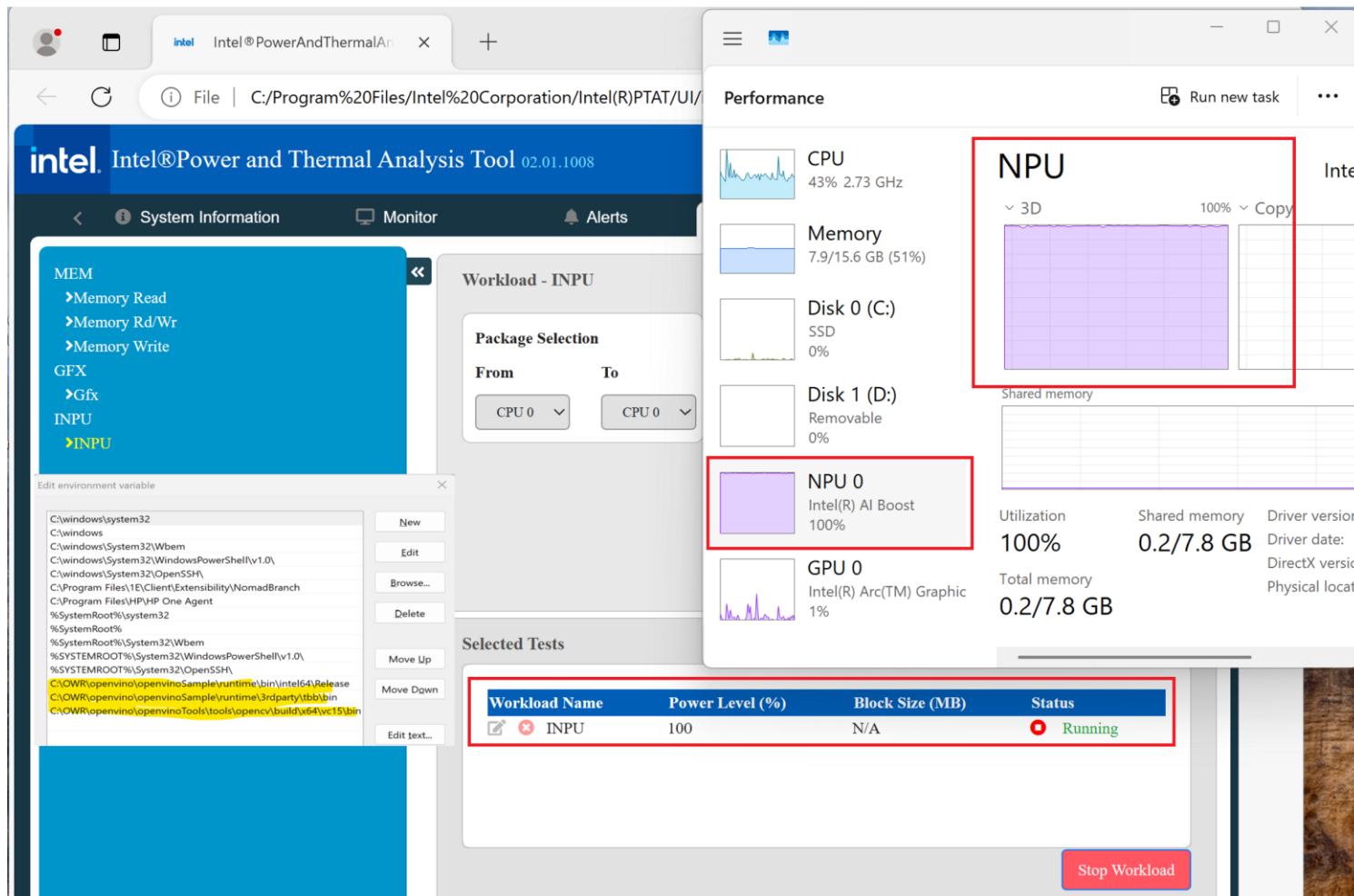
PTAT

- add the below three paths to the environment variables as follows (add to both systems and user path variables)



- copy(overwrite) that workload application "benchmark_app.exe" to the PTAT installed directory
C:\Program Files\Intel Corporation\Intel(R)PTAT\Vpu"

PTAT



MSFT WSE Tool

- Download [WSE Tool](#) (WseEnablingStatus.exe and WsePerformanceAssessmentTool.exe)
- WseEnablingStatus.exe
 - Checks the status of MEP installed on the system
 - System Name: device name
 - System OS Info: operating system info
 - Opt-In Camera Status: True/False
 - Opt-In Camera FriendlyName: friendly name of camera
 - Opt-In Camera Hardware ID: hardware id of camera
 - Opt-In Camera Driver: driver version of camera
 - Windows Studio Effects Camera: version number
 - Windows Studio Effects AudioBlur APO: version number (by platforms)
 - PerceptionCore.dll: version number [path]
 - You will see MEP opt-in camera and Windows Studio Effect version when you run this tool if you opt-in MEP successfully.

```
C:\WseEnablingStatus>WseEnablingStatus.exe
System Name: LAPTOP-CD63QRRN
System OS Info: Windows 10 Pro (26085.1)
Opt-In Camera Status: True
Opt-In Camera FriendlyName: 9MP Camera
Opt-In Camera Hardware ID: USB\VID_0408&PID_546E&REV_0006&MI_00
Opt-In Camera Driver: 10.0.26085.1
Windows Studio Effects Camera: 1.0.38.0
```

MSFT WSE Tool

- The tool will point out the problem if you don't opt-in MEP correctly.

```
F:\WseEnablingStatus>WseEnablingStatus.exe  
can not find 'Windows Studio Effects Camera' in device manager, extension .inf for MEP camera was not correctly deployed
```

```
F:\WseEnablingStatus>WseEnablingStatus.exe  
can not find Opt-in camera instance in registry, there is no 'FSMEnableMsEffects' key in registry
```

MSFT WSE Tool

- WsePerformanceAssessmentTool.exe

- Allows profiling the following metrics

- Camera Effects Performance:
 - FPS (FrameRate): the average frame per seconds of recorded clip
 - timeToFirstFrame: the time cost of successfully loading AI models and applying effects onto 1st camera frame
 - avgProcessingTime: the average processing time of dealing with one frame, in milliseconds
 - (numberOfFramesAbove33ms/numberOfProcessFrame): the ratio of processing time over 33 milliseconds
 - initTimeCameraApp: the time cost between launch camera app and applying effects onto 1st camera frame

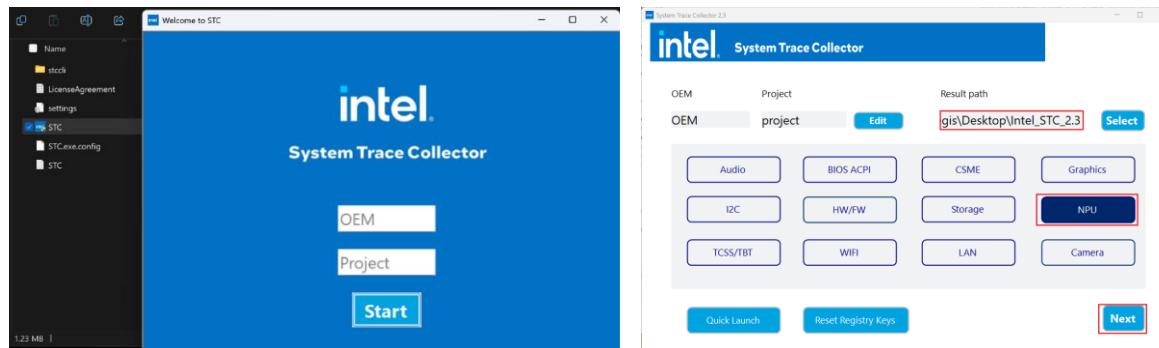
```
PerformanceAssessment [2 / 3] rounds
[Camera Effects Performance]
GenericError - PerceptionCore-v2.3.7
FPS (FrameRate): 24.36 frames/sec, (1920x1080, 00:00:59)
timeToFirstFrame: 6242 millisecs
avgProcessingTime: 21.69 millisecs
(numberOfFramesAbove33ms/numberOfProcessFrame): (0/1568)
InitTimeCameraApp: 7.048 secs

PerformanceAssessment [3 / 3] rounds
[Camera Effects Performance]
GenericError - PerceptionCore-v2.3.7
FPS (FrameRate): 25.34 frames/sec, (1920x1080, 00:00:59)
timeToFirstFrame: 6307 millisecs
avgProcessingTime: 21.29 millisecs
(numberOfFramesAbove33ms/numberOfProcessFrame): (0/1622)
InitTimeCameraApp: 7.116 secs
```

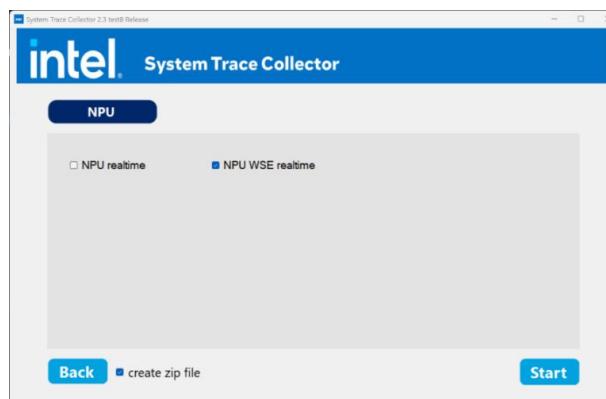
Logs/Debug

Intel STC v2.3

- To collect Camera and NPU UMD/KMD trace, Xperf log
- Download [Intel System Trace Collector \(STC\) Tool](#), launch it by STC.exe



- Select “NPU WSE realtime” (Camera and NPU UMD/KMD)



Intel STC v2.3

- You will see another console pop up.

```
Administrator: Windows PowerShell
Transcript started, output file is C:\users\Regis\AppData\Local\Temp\Regis_MediaTrace.log
Trace script log: C:\Users\Regis\AppData\Local\Temp\Regis_MediaTrace.log
Version: 1.1
Looking for the logging scenarios...
Gathering system information...
[Get-EnvironmentInformation] collecting environment information
Preparing local system...
Preparing target system...
Saving target system details...
Queue DxDiag to background job
Creating tracing scripts...
Starting tracing...

**** RUN YOUR SCENARIO NOW AND PRESS [ENTER] WHEN FINISHED ****
```

- Reproduce the issue and press ENTER when finished.

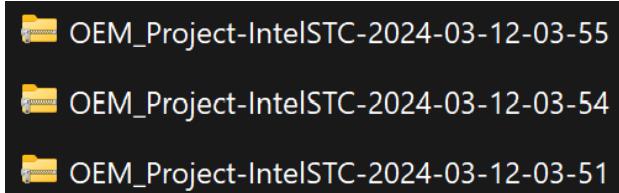
```
Administrator: Windows PowerShell
Transcript started, output file is C:\Users\MELAB\AppData\Local\Temp\TELAB_MediaTrace.log
Trace script log: C:\Users\MELAB\AppData\Local\Temp\TELAB_MediaTrace.log
Version: 1.1
Looking for the logging scenarios...
Gathering system information...
[Get-EnvironmentInformation] collecting environment information
Preparing local system...
Preparing target system...
Saving target system details...
Queue DxDiag to background job
Creating tracing scripts...
Starting tracing...

**** RUN YOUR SCENARIO NOW AND PRESS [ENTER] WHEN FINISHED ****

Stopping tracing and merging results...
Saving target system details...
Queue SetupAPI log to background job
Queue PnpUtil to background job
Queue WinHelloInfo to background job
Queue Winbio.evtx to background job
Queue MicrosoftTeamsLog to background job
Waiting for the background jobs to complete...
1 job(s) left.
```

Intel STC v2.3

- It will auto zip the log



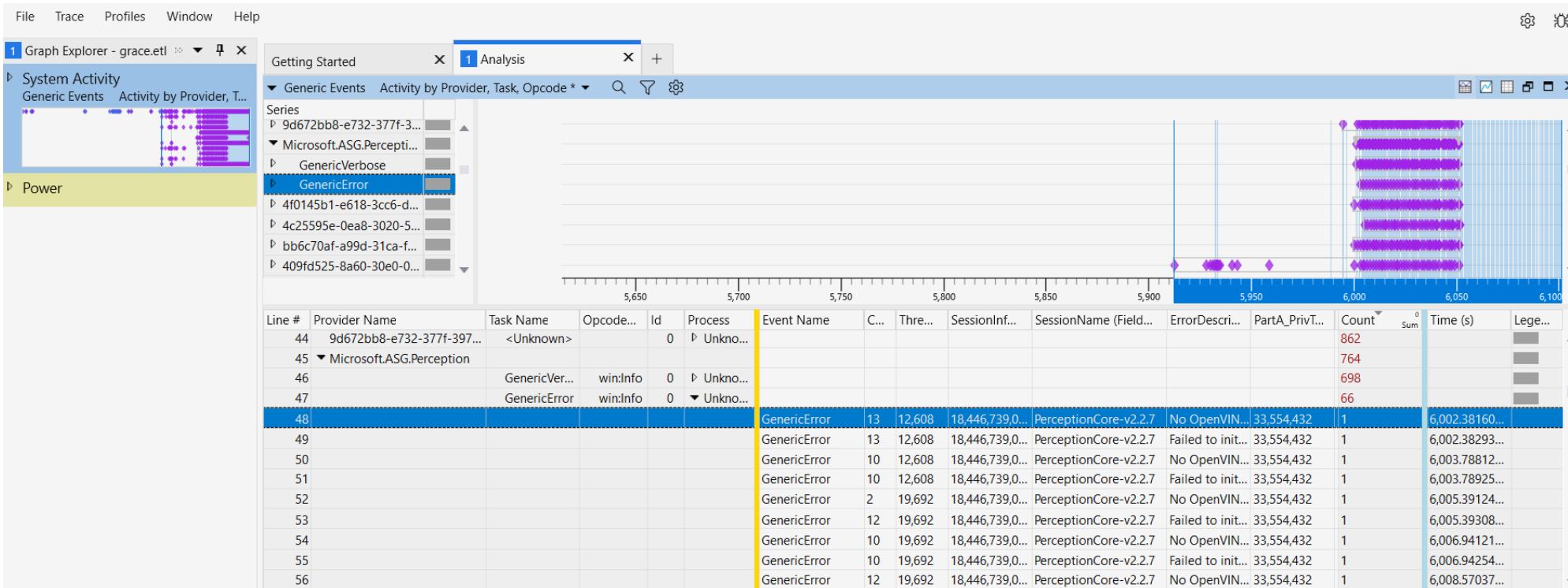
- Use Windows Performance Analyzer (WPA) tool in ADK to open Trace_Multimedia.etl

| Name | Date modified | Type | Size |
|------------------|-------------------|-----------------------|------------|
| Scripts | 3/12/2024 3:55 AM | File folder | |
| BuildInfo | 3/12/2024 3:55 AM | Text Document | 1 KB |
| dxdiag | 3/12/2024 3:55 AM | Text Document | 115 KB |
| pnpUtil.pnp | 3/12/2024 3:56 AM | PNP File | 6,627 KB |
| xxxxx_MediaTrace | 3/12/2024 3:56 AM | Text Document | 2 KB |
| setupapi.dev | 3/12/2024 3:14 AM | Text Document | 231 KB |
| Trace_Multimedia | 3/12/2024 3:56 AM | Windows Performanc... | 117,760 KB |
| winbio | 3/12/2024 3:56 AM | Event Log | 1,092 KB |
| WinHelloInfo | 3/12/2024 3:56 AM | Text Document | 4 KB |

Intel STC v2.3

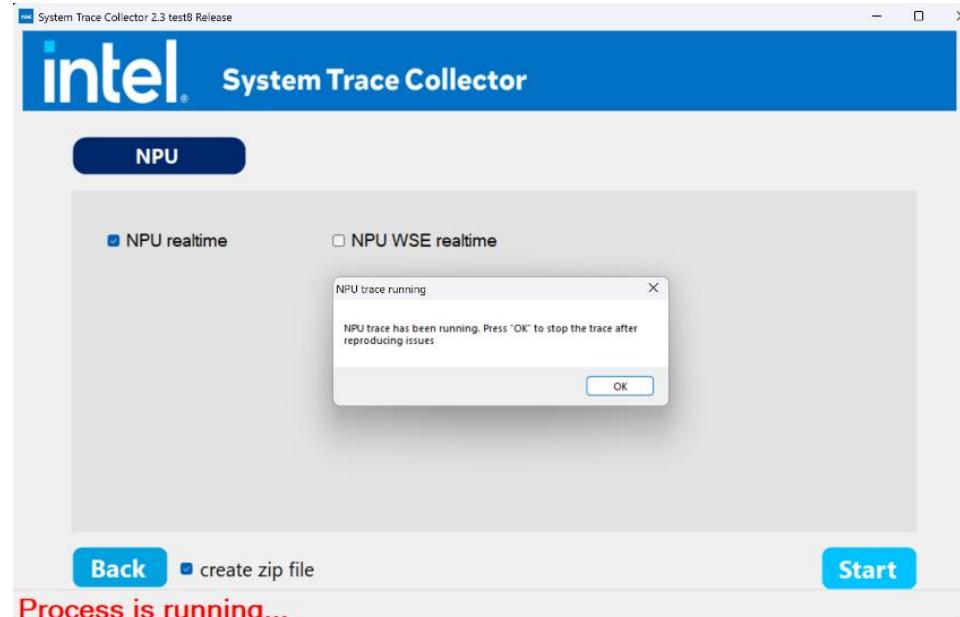
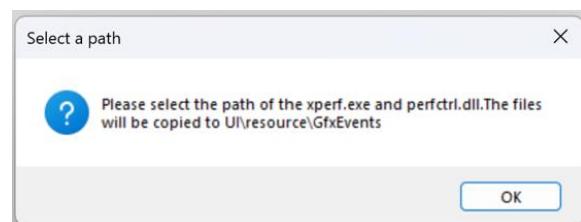
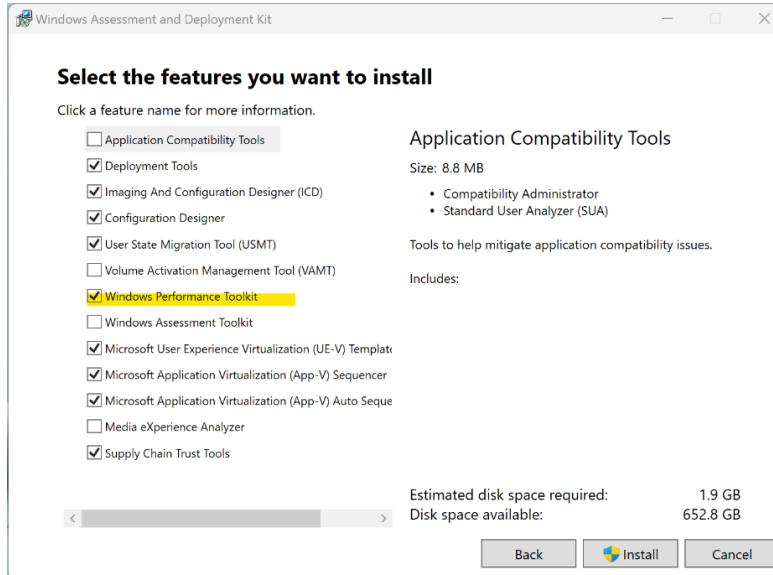
The following strings can be used to collect the MEP related logs.

- Microsoft.Windows.Capture.WindowsEffectCameraMediaSource
- Microsoft.ASG.Perception



Intel STC v2.3

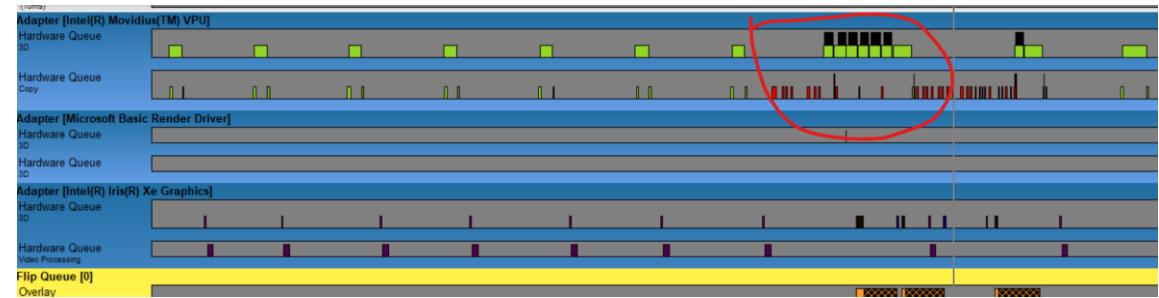
- Collect Xperf log by STC tool
- Pre-request: install [Windows ADK](#)



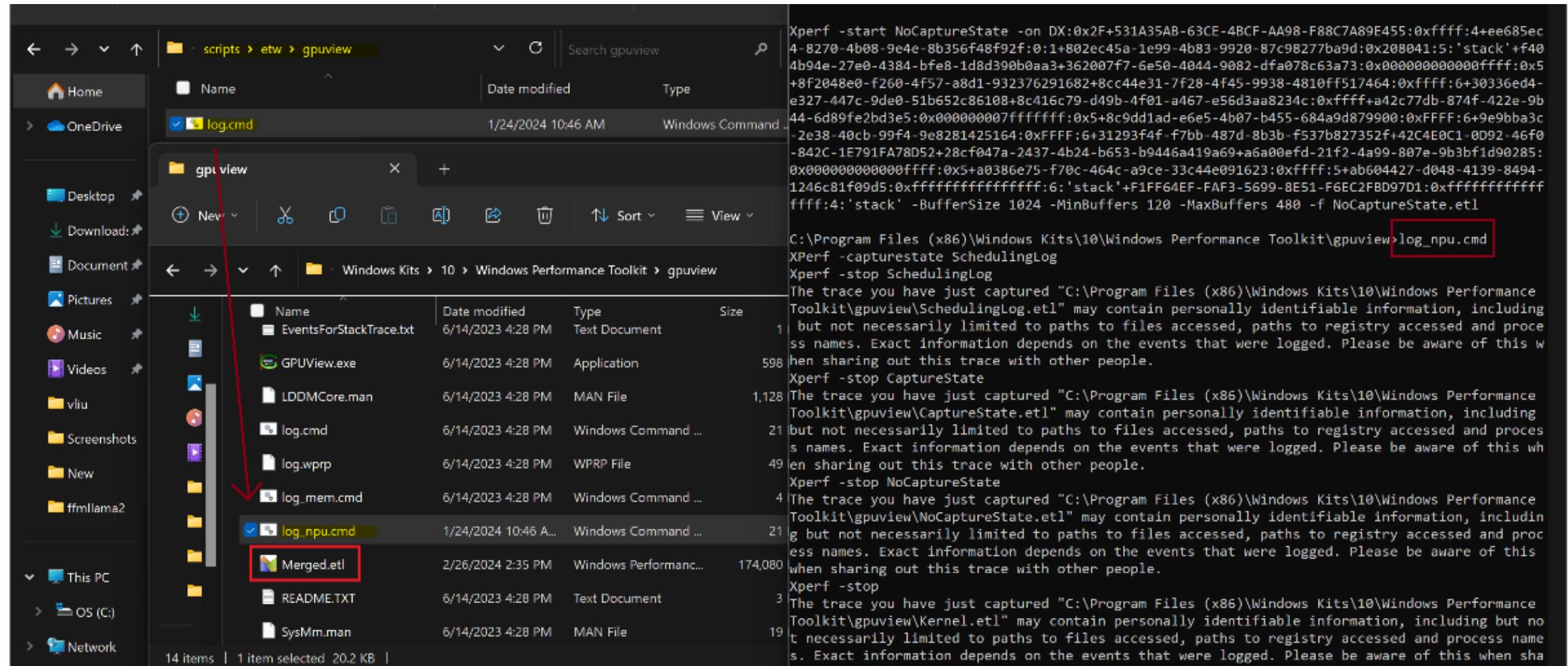
GPUView

GPUView is useful to debug performance issues. The adaptor NPU will record the inference event, click the block can get this inference details, including execution time, submission and complete time.

- Install Windows ADK
- In NPU driver release package, rename script/etw/gpuvview/log.cmd to log_npu.cmd and copy it to Windows Performance Toolkit\gpuvview folder.
- Execute cmd.exe as administrator and run log_npu.cmd.
- Reproduce the issue.
- When test complete, go back to the Admin cmd window and run log_npu.cmd again to generate the Merged.etl file



GPUView

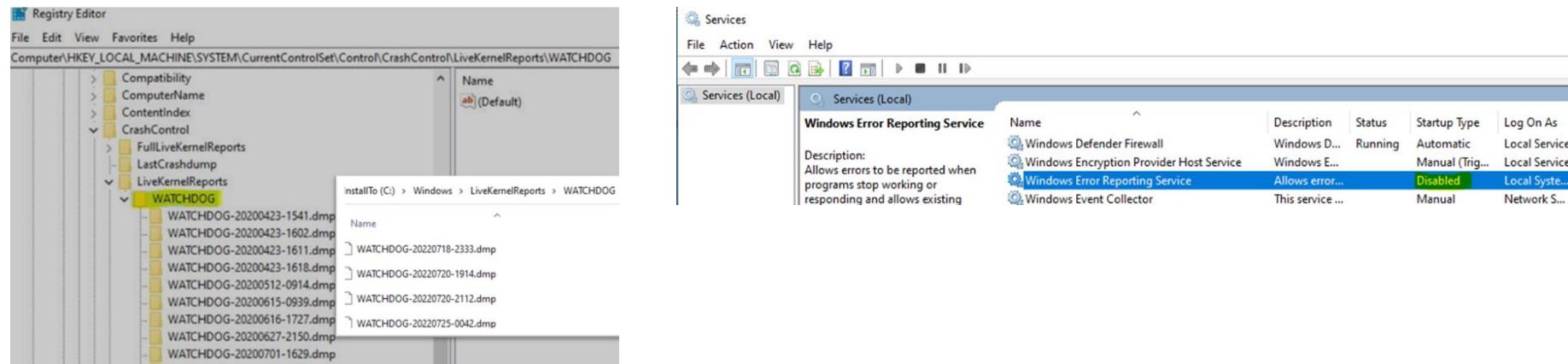


Windbg-Live kernel dump/BSOD memory dump

NPU driver also supports kernel live dump, we can use Windbg tool to check who causes the TDR.
The dump file is in C:\Windows\LiveKernelReports\WATCHDOG

If dump was not observed

- Make sure sub reg keys under Computer\HKEY_LOCAL_MACHINE\SYSTEM\CurrentControlSet\Control\CrashControl\LiveKernelReports\WATCHDOG are cleaned up, if it reaches 10, you won't see OS report files dump
- Disable the Windows Error Reporting Service



Windbg-Live kernel dump/BSOD memory dump

```
Loading User Symbols
Mini Kernel Dump does not contain unloaded driver list
For analysis of this file, run !analyze -v
watchdog!WdpDbgCaptureTriageDump+0xb7:
fffffb04`3356396d 488b4628    mov     qword ptr [rsi+28h] ds:002b:ffffbd0f`9499d038=??????
9: kd> !analyze -v
*****
*          *
*          Bugcheck Analysis
*          *
*****
VIDEO_ENGINE_TIMEOUT_DETECTED (141)
One of the display engines failed to respond in timely fashion.
(This code can never be used for a real BugCheck.)
Arguments:
Arg1: fffffbd0f94d6c010, Optional pointer to internal TDR recovery context (TDR_RECOVERY_CONTEXT).
Arg2: fffff80457c03330, The pointer into responsible device driver module (e.g. owner tag).
Arg3: 0000000000000000, The secondary driver specific bucketing key.
Arg4: 00000000000309c, Optional internal context dependent data.
```

```
FILE_IN_CAB: WATCHDOG-20230315-0332.dmp
DUMP_FILE_ATTRIBUTES: 0x18
Kernel Generated Triage Dump
Live Generated Dump
BUGCHECK_CODE: 141
BUGCHECK_P1: fffffbd0f94d6c010
BUGCHECK_P2: fffff80457c03330
BUGCHECK_P3: 0
BUGCHECK_P4: 309c
TAG_NOT_DEFINED_202b: *** Unknown TAG in analysis list 202b
VIDEO_TDR_CONTEXT: dt dxgkrnl!_TDR_RECOVERY_CONTEXT fffffbd0f94d6c010
Symbol dxgkrnl!_TDR_RECOVERY_CONTEXT not found.
PROCESS_OBJECT: 000000000000309c
PROCESS_NAME: System
STACK_TEXT:
fffffdc8d'6f9a7550 fffff804`33578f24 : fffffbd0f`8768b030 fffffbd0f`8768b030 fffffbd0f`9499d010 fffffbd0f`94d6c010 : watchdog!WdpDbgCaptureTriageDump+0xb7
fffffdc8d'6f9a75c0 fffff804`333dbd0 : fffffbd0f`94d6c010 fffffdc8d`6f9a7780 00000000`00000000 fffff804`1ba54e9 : watchdog!WdpDbgReportRecreate+0xd4
fffffdc8d'6f9a7620 fffff804`1ba2bf5b : fffffbd0f`00000000 fffffbd0f`90b36bd0 fffffbd0f`94d6c010 fffffbd0f`90b3c000 : dxgkrnl!TdrUpdateDbgReport+0x11d
fffffdc8d'6f9a7680 fffff804`1ba5f589 : fffffbd0f`90b3c000 00000000`00000000 fffffbd0f`90b36000 fffffbd0f`90b3c001 : dxgms2!VidSchiResetEngine+0x70f
fffffdc8d'6f9a7830 fffff804`1ba512b : fffffbd0f`90b36000 00000000`00000000 00000000`00000000 00000000`00000000 : dxgms2!VidSchiResetEngines+0xb1
fffffdc8d'6f9a7880 fffff804`1ba97c9f : 00000000`00000000 00000000`00000000 00000000`000029e0 00000000`00989680 : dxgms2!VidSchiCheckHwProgress+0x2d44b
fffffdc8d'6f9a7900 fffff804`1b9fb689 : 00000000`00000000 fffffbd0f`90b36000 fffffdc8d`6f9a7a39 00000000`00000000 : dxgms2!VidSchiWaitForSchedulerEvents+0x37f
fffffdc8d'6f9a79d0 fffff804`1baa3a95 : fffffbd0f`8acbbaa0 fffffbd0f`90b36000 fffffbd0f`8acbbaa0 fffffbd0f`961c6010 : dxgms2!VidSchiScheduleCommandToRun+0x309
fffffdc8d'6f9a7aa0 fffff804`1baa3a8a : 00000000`00000000 fffff804`1baa3940 fffffbd0f`90b36000 00000000`00000080 : dxgms2!VidSchiRun_PriorityTable+0x35
fffffdc8d'6f9a7af0 fffff804`17c0f5b7 : fffffbd0f`90b39100 fffff804`00000001 fffffbd0f`90b36000 006fe47f`b19bd0ff : dxgms2!VidSchiWorkerThread+0xca
fffffdc8d'6f9a7b30 fffff804`17e2e364 : fffffe081`1d9a1180 fffffbd0f`90b39100 fffff804`17c0f560 97f29678`aa1603e0 : nt!PspSystemThreadStartup+0x57
fffffdc8d'6f9a7b80 00000000`00000000 : fffffdc8d`6f9a1000 00000000`00000000 00000000`00000000 : nt!KiStartSystemThread+0x34
SYMBOL_NAME: ivdkmd+3330
MODULE_NAME: ivdkmd
IMAGE_NAME: ivdkmd.sys
```

Windbg-Live kernel dump/BSOD memory dump

- If the system shows blackscreen or blue screen (BSOD), see if any memory dump created under C:\Windows\MEMORY.DMP.
- Using Windbg tool to check if NPU driver npu_kmd.sys or ivdkmd.sys in Windbg !analyze -v

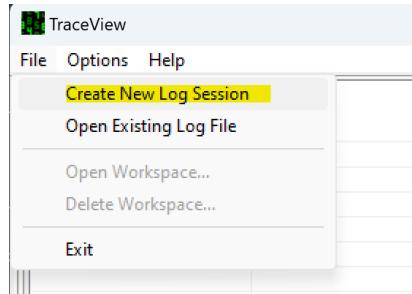
```
2: kd> !analyze -v
*****
*          Bugcheck Analysis
*
*****
VIDEO_TDR_FAILURE (116)
Attempt to reset the display driver and recover from timeout failed.
Arguments:
Arg1: ffff810fa50c6400, Optional pointer to internal TDR recovery context (TDR_RECOVERY_CONTEXT).
Arg2: fffff80583073000, The pointer into responsible device driver module (e.g. owner tag).
Arg3: 0000000000000000, Optional error code (NTSTATUS) of the last failed operation.
Arg4: 000000000000000d, Optional internal context dependent data.

STACK_TEXT:
fffffa81'00e776f8 fffff805'6e2a3e8d : 00000000 00000016 ffff810f`5dc6400 fffff805'83073000 00000000 00000000 : nt!KeBugCheckEx
fffffa81'00e776d0 fffff805'6e2a3e8d : fffff805'83073000 00000010`5dc6400 fffff810f`00e77630 00000000 00000000 : nt!KeBugCheckOnTimeout+0x1fe
fffffa81'00e77740 fffff805'64f45189 : 00000000 00000004 fffff810f`a3e6a000 01d81c05`4d4fffd8 fffff805'64f45189 : nt!KeBugCheckOnTimeout+0x19d
fffffa81'00e77770 fffff805'64f45189 : 00000000 00000164 fffff810f`a3e6a000 01d81c05`4d4fffd8 fffff805'64f45189 : nt!KeBugCheckOnTimeout+0x19d
fffffa81'00e77780 fffff805'6505826c : fffff810f`a3e6a000 00000000 00000001 ffff810f`a3e6a000 00000000 00000000 : nt!KeBugCheckOnTimeout+0x19d
fffffa81'00e77900 fffff805'64f7b949 : 00000000 00000000 ffff810f`a3e6a000 fffff810f`00e77a39 00000000 00000000 : nt!KeBugCheckOnTimeout+0x19d
fffffa81'00e77940 fffff805'6501e335 : 00000000 00000000 ffff810f`a3e6a000 fffff810f`a1c0b010 fffff810f`a1c11010 : nt!KeBugCheckOnTimeout+0x19d
fffffa81'00e77980 fffff805'6501e335 : fffff810f`a1c0b010 fffff810f`a3e6a000 fffff810f`a1c11010 00000000 00000000 : nt!KeBugCheckOnTimeout+0x19d
fffffa81'00e77a40 fffff805'6501e335 : fffff810f`a1c0b010 fffff810f`a3e6a000 fffff810f`a1c11010 00000000 00000000 : nt!KeBugCheckOnTimeout+0x19d
fffffa81'00e77a80 fffff805'6501e335 : fffff810f`a1c0b010 fffff810f`a3e6a000 fffff810f`a1c11010 00000000 00000000 : nt!KeBugCheckOnTimeout+0x19d
fffffa81'00e77b30 fffff805'6804765 : fffff810f`a3e020c0 fffff805'6501e335 00000001 ffff810f`a3e6a000 006fe47f`b19bbdff : nt!PspSystemThreadStartup+0x55
fffffa81'00e77b30 fffff805'68219b4 : fffff810f`a3e020c0 fffff805'680478a0 48c03302`ebc00000 : nt!PspSystemThreadStartup+0x55
fffffa81'00e77b80 00000000 00000000 : fffff810f`00e78000 fffffaf81'00e71000 00000000 00000000 00000000 : nt!KiStartSystemThread+0x34

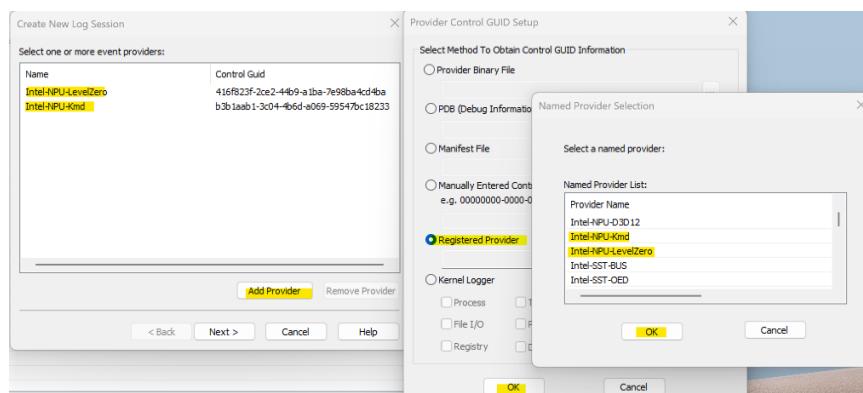
SYMBOL_NAME: ivdkmd+3000
MODULE_NAME: ivdkmd
IMAGE_NAME: ivdkmd.sys
IMAGE_VERSION: 29.21.10.257
STACK_COMMAND: .cxr; .ecxr ; kb
FAILURE_BUCKET_ID: 0x116_IMAGE_ivdkmd.sys
```

Trace View

- Install Windows Driver Kit (WDK)
 - The traceview.exe can be found in C:\Program Files (x86)\Windows Kits\10\Tools\10.0.22621.0\x64\traceview.exe
- Right click traceview.exe, select “Run as administrator”. “File” -> “Create New Log Session”.

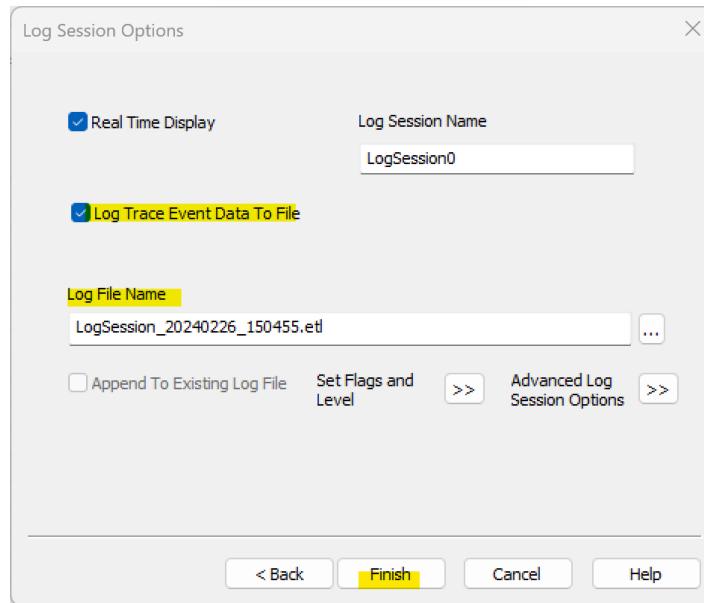


- In Registered Provider, add Intel-NPU-Kmd and Intel-NPU-LevelZero



Trace View

- Click Next and check “Log Trace Event Data To File” and then “Finish” to start the log session.
- Reproduce the issue
- Selecting “File”> “Exit” to stop the log session



| TraceView | | | | | | | | | | | | | | | |
|---|---------|------------------------|-----------|------|-----------|--------------|---------------------------------|--|--|--|--|--|--|--|--|
| File | | Options | | Help | | | | | | | | | | | |
| Create New Log Session | | Open Existing Log File | | | | | | | | | | | | | |
| Event Count | | Events Lost | | | | | | | | | | | | | |
| 16424 | | 0 | | | | | | | | | | | | | |
| Buffers Read Flags Max Buf Min Buf Level KD Filter Ignore TraceView Max Trace Records Log File Name | | | | | | | | | | | | | | | |
| 48 * 21 4 * FALSE FALSE 65536 LogSession_2... | | | | | | | | | | | | | | | |
| Exit | | | | | | | | | | | | | | | |
| Msg# | Name | Process ID | Thread ID | CPU# | Sequence# | System Time | Message | | | | | | | | |
| 000... | Inte... | 9536 | 16724 | 17 | 0 | 02\26\202... | SetPowerComponentActive - Start | | | | | | | | |
| 000... | Inte... | 9536 | 16724 | 17 | 0 | 02\26\202... | SetPowerComponentActive - End | | | | | | | | |
| 000... | Inte... | 4 | 1252 | 10 | 0 | 02\26\202... | BuildPagingBuffer - Start | | | | | | | | |
| 000... | Inte... | 4 | 1252 | 10 | 0 | 02\26\202... | BuildPagingBuffer - End | | | | | | | | |
| 000... | Inte... | 4 | 1252 | 10 | 0 | 02\26\202... | BuildPagingBuffer - Start | | | | | | | | |
| 000... | Inte... | 4 | 1252 | 10 | 0 | 02\26\202... | BuildPagingBuffer - End | | | | | | | | |
| 000... | Inte... | 4 | 1252 | 10 | 0 | 02\26\202... | BuildPagingBuffer - Start | | | | | | | | |
| 000... | Inte... | 4 | 1252 | 10 | 0 | 02\26\202... | SetPowerComponentActive - Start | | | | | | | | |
| 000... | Inte... | 4 | 1252 | 10 | 0 | 02\26\202... | SetPowerComponentActive - End | | | | | | | | |
| 000... | Inte... | 0 | 0 | 0 | 0 | 02\26\202... | InterruptRoutine - Start | | | | | | | | |
| 000... | Inte... | 0 | 0 | 0 | 0 | 02\26\202... | InterruptRoutine - End | | | | | | | | |
| 000... | Inte... | 0 | 0 | 0 | 0 | 02\26\202... | DpcRoutine - Start | | | | | | | | |
| 000... | Inte... | 0 | 0 | 0 | 0 | 02\26\202... | DpcRoutine - End | | | | | | | | |

| This PC > OS (C) > Program Files (x86) > Windows Kits > 10 > Tools > 10.0.22621.0 > x64 | | | | |
|---|--------------------------------|-------------------|------------------------|-----------|
| | Name | Date modified | Type | Size |
| | LogSession_20240226_152108.etl | 2/26/2024 3:27 PM | Windows Performance... | 17,200 KB |
| | wsdbit_server.exe | 9/30/2023 2:19 AM | Application | 66 KB |
| | sensordiagnostictool.exe | 9/30/2023 2:19 AM | Application | 274 KB |

WPR (Windows Performance Recorder)

Next Gen AI Experiences - NPU ETL recording and analysis

- wpr.exe -start NeuralProcessing -filemode
 - repro issue
 - wpr -stop trace.etl
 - can be combined with other built-in profiles or your own tracing profile(.WPRP file)
 - ex: wpr -start NeuralProcessing -start CPU -start file.wprp!myprofile

AITrace.wprp

[4/15/2024 Update] Version 1.1 of this package comes with a AITrace.wprp custom tracing profile that contains some ONNX/DirectML and other core AI framework ETW instrumentation. A few providers from inbox Windows AI experiences are also included (Paint, Live Caption, Recall, etc.). It can be used on top of the NeuralProcessing flag to obtain additional performance markers.

NPU ETL recording and analysis

- ADK 26063+ [Download the Windows ADK 10.1.26100.1 \(May 2024\)](#)
- It shows each process consuming the resource, the callstack dispatching the work to the NPU, and the amount of time spent consuming the resource(in absolute time and relative % based on the time interval focused on)

Select the features you want to install

Click a feature name for more information.

- Application Compatibility Tools
- Deployment Tools
- Imaging And Configuration Designer (ICD)
- Configuration Designer
- User State Migration Tool (USMT)
- Volume Activation Management Tool (VAMT)
- Windows Performance Toolkit
- Windows Assessment Toolkit
- Microsoft User Experience Virtualization (UE-V) Template Generator

Supply Chain Trust Tools

Size: 201.4 MB

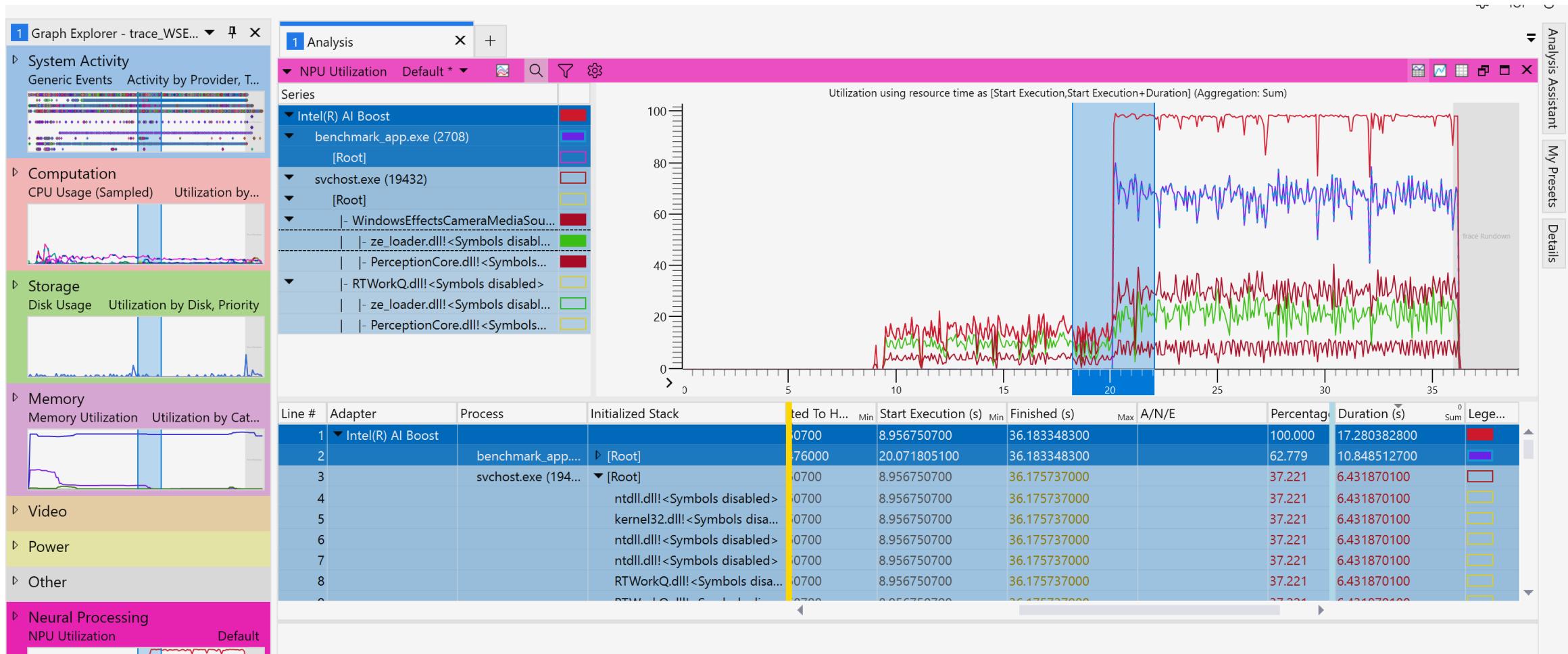
- SBOM Tool
- CoseSign Tool

Tools for creating SBOMs and managing trust.

Includes:



NPU ETL recording and analysis



PTAT

For NPU performance issues, we can also get NPU frequency data by [PTAT](#) tool

The screenshot shows the Intel® Power and Thermal Analysis Tool (PTAT) interface. The top navigation bar includes tabs for System Information, Monitor (which is selected), Alerts, Workload, Control, Scripts, and Grid View. It also displays the connection status "Connected to : DESKTOP-6REFH23". Below the navigation bar, there are two main sections: a left sidebar with monitoring categories and a central data grid.

Left Sidebar (Monitoring Categories):

- Monitor All
- Battery
- CPU
- CState
- HGS
- Memory
- PCH
- Topology
- iNPU

Central Data Grid (Grid View):

The central grid displays data for the "iNPU" component, specifically the "Integrated NPU". The columns are "Type" and "Value". The data includes:

| Type | Value |
|------------------------|-------|
| NPU Device ID | NA |
| Tile Count | NA |
| DPU Count | NA |
| Shave Count | NA |
| DMA Count | NA |
| CMX size per tile(MB) | NA |
| NPU Freq(MHz) | NA |
| Dpu_Clk(MHz) | NA |
| NPU_Clk(MHz) | NA |
| Minimum Frequency(MHz) | NA |
| Maximum Frequency(MHz) | NA |
| Turbo Capability | NA |

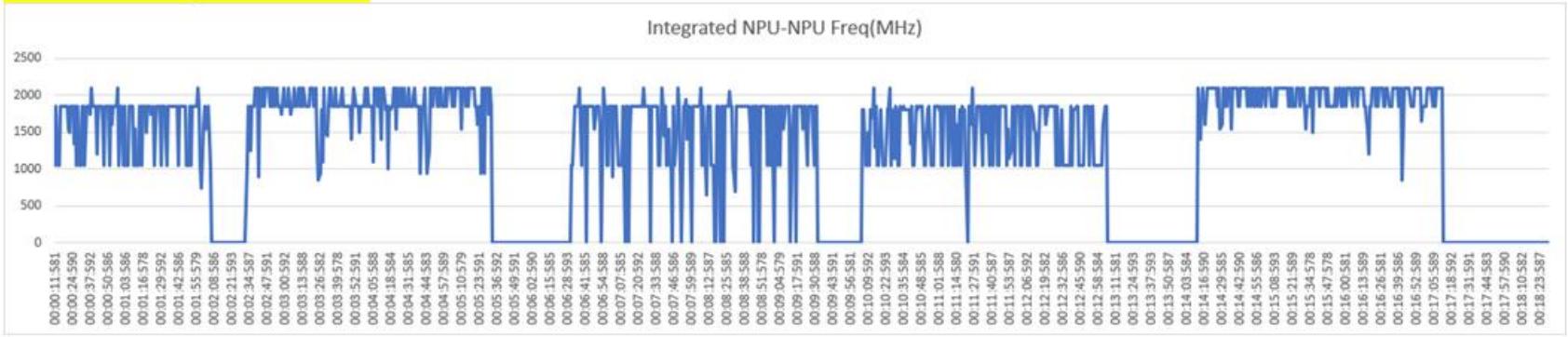
Right Panel (CPU Topology):

This panel lists various CPU topology parameters, each with a "Type" and "Value" column. Most values are marked as "NA".

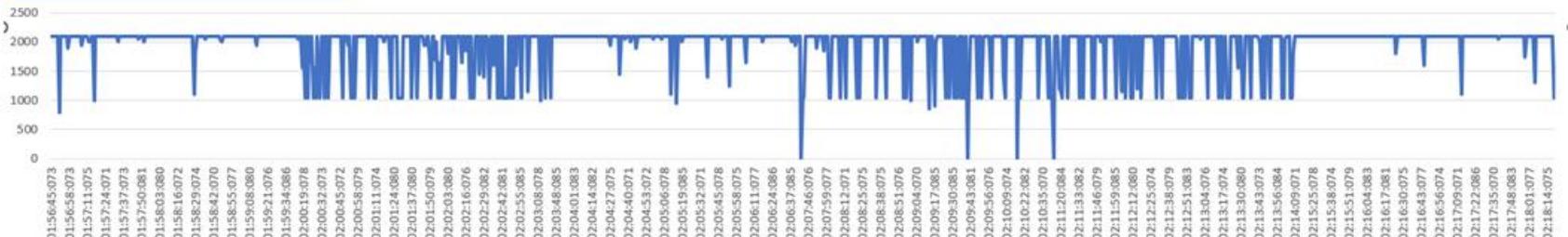
| Type | Value |
|------------------------------|-------|
| Number of Physical Processor | NA |
| Number of Cores | NA |
| Number of logical Cores | NA |
| HT Capability | NA |
| Physical Core-0 | NA |
| Physical Core-1 | NA |
| Physical Core-2 | NA |
| Physical Core-3 | NA |
| Physical Core-4 | NA |
| Physical Core-5 | NA |
| Physical Core-6 | NA |
| Physical Core-7 | NA |
| Physical Core-8 | NA |
| Physical Core-9 | NA |

PTAT

HP Ernesto: Procyon Score is 461



HP Willie: Procyon Score is 511



| Integrated NPU-NPU Freq(MHz) | Integrated NPU-Dpu_Clk(MHz) | Integrated NPU-NPU_Clk(MHz) |
|------------------------------|-----------------------------|-----------------------------|
| 950 | 633 | 475 |
| 950 | 633 | 475 |
| 950 | 633 | 475 |
| 950 | 633 | 475 |
| 950 | 633 | 475 |
| 1450 | 966 | 475 |
| 950 | 633 | 475 |
| 950 | 633 | 475 |
| 950 | 633 | 475 |
| 1550 | 1033 | 475 |
| 950 | 633 | 475 |
| 1000 | 666 | 500 |
| 950 | 633 | 475 |
| 950 | 633 | 475 |
| 950 | 633 | 475 |
| 950 | 633 | 475 |
| 950 | 633 | 475 |
| 950 | 633 | 475 |
| 950 | 633 | 475 |
| 950 | 633 | 475 |
| 950 | 633 | 475 |
| 950 | 633 | 475 |
| 950 | 633 | 475 |
| 950 | 633 | 475 |
| 950 | 633 | 475 |
| 950 | 633 | 475 |
| 950 | 633 | 475 |
| 950 | 633 | 475 |
| 950 | 633 | 475 |
| 950 | 633 | 475 |
| 1050 | 700 | 525 |
| 950 | 633 | 475 |
| 950 | 633 | 475 |

Reference

ARL NPU Introduction and Features (video)

2023 2H TPnP Customer Tour - MTL-TnP-O025 - MTL Resnet50 BKM (video)

MTL-AI-D006MTL iVPU Validation and Debug Training (video)

Meteor Lake VPU MEP Enablement and Validation Testcases in OEM Platform

Debug NPU/MEP software stack debug guide

The Intel logo is displayed in white against a solid blue background. The word "intel" is written in a lowercase, sans-serif font. A small, solid blue square is positioned above the letter "i". The letter "i" has a vertical stroke extending upwards from its top loop. The letter "t" has a vertical stroke extending downwards from its top loop. The letters "n", "e", and "l" are standard lowercase forms.