# RST POC Driver - DSM Hints

This POC driver is used to evaluate new SSDs which supports DSM Hints feature.

So far, focus on two use cases, refer to below:



## Evaluation Requirements:

1. Intel platform (ex: ARL-H system) - VMD controller enabled and SSDs under VMD
2. RST POC Driver 20.2 and Tools - Link
3. POC SSDs supporting DSM Hints
4. RSTCLI tool in RST 20.2.6.1025 kit

## Evaluate Use Case#1 (LLMs files - load/unload)

- Enable NVMe DSM Hints, use below two commands (if NVMe device ID is 2-0-0)

  NvmePassthroughApp.exe --scsi 0 --path 2 --target 0 --lun 0 configureDsm --enableNvmeHinting 1 --userModeHinting 1 --pageFileHinting 0 --readHinting 1 --writeHinting 0
  NvmePassthroughApp.exe --scsi 0 --path 2 --target 0 --lun 0 addDsmClassification --kind  2 --path \path_to_model_directory\

- Intel OpenVINO AI app for benchmark:

  **openvino.genai/samples/cpp/text_generation/benchmark_genai.cpp at master · openvinotoolkit/openvino.genai · GitHub**

  **or, use the 2025.4.1.0_x86_64 installer in storage.openvinotoolkit.org**
  **command options:**
  **-m <Llama2 7B INT4 OV model> -d GPU -p "The Sky is blue because"  --nw 0 -n 1 --mt 20**

  benchmark_genai.exe Help command
  Usage:
    benchmark_vanilla_genai [OPTION...]

```
-m, --model arg     Path to model and tokenizers base directory
-p, --prompt arg    Prompt (default: "")
   --pf arg         Read prompt from file
   --nw arg         Number of warmup iterations (default: 1)
-n, --num_iter arg  Number of iterations (default: 3)
   --mt arg         Maximal number of new tokens (default: 20)
-d, --device arg    device (default: CPU)
-h, --help          Print usage
```

for example:

benchmark_genai.exe -m C:\dev\pyworkspace\GenAI\models\open_llama_7b_v2-int4-ov -d GPU -p "The Sky is blue because" --nw 0 -n 1 --mt 20
OpenVINO Runtime
    Version : 2025.4.1
    Build   : 2025.4.1-20426-82bbf0292c5-releases/2025/4

Prompt token size:6
Output token size:20
Load time: 11660.00 ms
Generate time: 2532.16  0.00 ms
Tokenization time: 0.31  0.00 ms
Detokenization time: 0.62  0.00 ms
TTFT: 181.35  0.00 ms
TPOT: 123.68  12.68 ms/token
Throughput: 8.09  0.83 tokens/s

- Llama2 7B INT4 OV model:

  https://huggingface.co/OpenVINO/open_llama_7b_v2-int4-ov

- Need to modify the code of "benchmark_genai" sample in order to have a compiled model cache in a directory
  benchmark_genai -m .\models\open_llama_7b_v2-int4-ov -d GPU -p "The Sky is blue because" --nw 0 -n 1 --mt 20 --cache_dir ".ccache"

Compiled Cache Dir: compiled_cache
OpenVINO Runtime
   Version : 2025.4.1
   Build   : 2025.4.1-20426-82bbf0292c5-releases/2025/4

Using CACHE_DIR: .ccache
Prompt token size:6
Output token size:20
Load time: 5860.00 ms
Generate time: 1850.92  0.00 ms
Tokenization time: 0.53  0.00 ms
Detokenization time: 0.51  0.00 ms
TTFT: 131.27  0.00 ms
TPOT: 90.47  19.17 ms/token
Throughput: 11.05  2.34 tokens/s

## Some Discussions/ARs:

1. how to use the driver on ARL-H platform?
2. how to utilize Intel OpenVINO AI app to verify TTFT
3. POC SSDs
4. checkpoint schedule

# RST Driver change history:

v3: RST 20.2.0.8335
Fix of ioctl

v2:
change AccessLatency from 2 to 3 for files, processes, directories

v1:
for files, processes, directories: AccessFrequency = 3, AccessLatency = 3

ⓘ

## Related articles

- [RST POC Driver - DSM Hints](#)