

CS6886W - System Engineering for Deep Learning – Assignment2

Performance Benchmarking, Roofline Modeling, and Scaling Study

Sanchari Roy[CS24M532]

GitHub repository : <https://github.com/cs24m532/Deep-Learning>

Overview:

This report presents a comprehensive performance analysis of the **GPT-2 Medium** language model using the open-source **llama.cpp** inference framework. The study evaluates the model under multiple execution configurations, beginning with naive single-threaded CPU runs and progressing toward highly optimized builds that incorporate SIMD acceleration and Intel® MKL-based matrix multiplication. Each configuration is benchmarked to assess its computational efficiency, scalability, and hardware utilization.

The primary objective of this work is to understand how **software optimizations**, **compiler flags**, and **hardware parallelism** affect the inference performance of large language models. By leveraging the **Roofline performance model**, the study quantifies key factors such as **operational intensity (OI)**, **memory bandwidth usage**, and **achieved FLOP/s**, enabling a deeper examination of the system's memory-bound versus compute-bound behavior.

Experimental results reveal clear trade-offs between execution regimes: naive builds suffer from significant memory stalls, while optimized and multi-threaded builds demonstrate notable speedups—though ultimately constrained by the system's memory bandwidth ceiling. Through this systematic evaluation, the report provides valuable insights into **hardware-aware optimization**, **scaling behavior**, and **performance tuning techniques** essential for efficient deep learning inference workloads on modern CPU architectures.

Since I am not able to find any single machine where I can get all the counters I have to switch between windows and Linux machines

Task 1 – Installing llama.cpp from Source

Steps to Execute:

`git clone https://github.com/ggml-org/llama.cpp`

-----Clones the llama.cpp repository.

`cd llama.cpp`

-----Moves into the cloned project directory

`cmake -B build`

----- CMake configures the build system and generates necessary makefiles.

`cmake --build build --config Release -j 8`

-----The build uses 8 parallel jobs for faster compilation.Verification:

A successful build produces the binary build/bin/llama-bench.

Ran `./build/bin/llama-bench --help` confirms installation.

Snapshots:

```
C:\Users\sanchroy>git version
git version 2.51.2.windows.1

C:\Users\sanchroy>cd Desktop

C:\Users\sanchroy\Desktop>cd MTECH

C:\Users\sanchroy\Desktop\MTECH>cd "deep learning"

C:\Users\sanchroy\Desktop\MTECH\deep learning>git clone https://github.com/ggml-org/llama.cpp
Cloning into 'llama.cpp'...
remote: Enumerating objects: 66485, done.
remote: Counting objects: 100% (118/118), done.
remote: Compressing objects: 100% (42/42), done.
remote: Total 66485 (delta 96), reused 76 (delta 76), pack-reused 66367 (from 4)
Receiving objects: 100% (66485/66485), 185.80 MiB | 6.24 MiB/s, done.
Resolving deltas: 100% (48290/48290), done.
Updating files: 100% (1845/1845), done.

C:\Users\sanchroy\Desktop\MTECH\deep learning>
```

```
C:\Users\sanchroy\Desktop\MTECH\deep learning\llama.cpp>
C:\Users\sanchroy\Desktop\MTECH\deep learning\llama.cpp>cmake -B build -DLLAMA_CURL=OFF
CMAKE_BUILD_TYPE=Debug
-- Warning: cache not found - consider installing it for faster compilation or disable this warning with GGML_CCACHE=OFF
-- CMAKE_SYSTEM_PROCESSOR: AMD64
-- CMAKE_GENERATOR_PLATFORM:
-- GGML_SYSTEM_ARCH: x86
-- Including CPU backend
-- x86 detected
-- Adding CPU backend variant ggml-cpu: /arch:AVX GGML_AVX
-- ggml version: 0.9.4
-- ggml_commit: 1ae74882f
-- Configuring done (1.6s)
-- Generating done (2.1s)
-- Build files have been written to: C:/Users/sanchroy/Desktop/MTECH/deep learning/llama.cpp/build

C:\Users\sanchroy\Desktop\MTECH\deep learning\llama.cpp>cmake --build build --config Release -j 8
Warning: NMake does not support parallel builds. Ignoring parallel build command line option.
[  0%] Building C object ggml/src/CMakeFiles/ggml-base.dir/ggml.c.obj
ggml.c
[  1%] Building CXX object ggml/src/CMakeFiles/ggml-base.dir/ggml.cpp.obj
ggml.cpp
[  1%] Building C object ggml/src/CMakeFiles/ggml-base.dir/ggml-alloc.c.obj
ggml-alloc.c
[  1%] Building CXX object ggml/src/CMakeFiles/ggml-base.dir/ggml-backend.cpp.obj
ggml-backend.cpp
[  1%] Building CXX object ggml/src/CMakeFiles/ggml-base.dir/ggml-opt.cpp.obj
ggml-opt.cpp
[  2%] Building CXX object ggml/src/CMakeFiles/ggml-base.dir/ggml-threading.cpp.obj
ggml-threading.cpp
[  2%] Building C object ggml/src/CMakeFiles/ggml-base.dir/ggml-quants.c.obj
ggml-quants.c
```

```
cvector-generator.cpp
C:\Users\sanchroy\Desktop\MTECH\deep learning\llama.cpp\common\common.h(259): warning C4305: 'initializing': truncation from 'double' to 'float'
C:\Users\sanchroy\Desktop\MTECH\deep learning\llama.cpp\common\common.h(269): warning C4244: 'argument': conversion from 'const unsigned int' to 'float', possible loss of data
C:\Users\sanchroy\Desktop\MTECH\deep learning\llama.cpp\tools\cvector-generator\pca.hpp(38): warning C4305: 'initializing': truncation from 'double' to 'float'
C:\Users\sanchroy\Desktop\MTECH\deep learning\llama.cpp\tools\cvector-generator\pca.hpp(109): warning C4244: 'argument': conversion from 'int64_t' to 'const unsigned int', possible loss of data
C:\Users\sanchroy\Desktop\MTECH\deep learning\llama.cpp\tools\cvector-generator\cvector-generator.cpp(101): warning C4244: 'argument': conversion from 'float' to 'const unsigned int', possible loss of data
C:\Users\sanchroy\Desktop\MTECH\deep learning\llama.cpp\tools\cvector-generator\cvector-generator.cpp(102): warning C4244: 'argument': conversion from 'float' to 'const unsigned int', possible loss of data
C:\Users\sanchroy\Desktop\MTECH\deep learning\llama.cpp\tools\cvector-generator\cvector-generator.cpp(103): warning C4244: 'argument': conversion from 'float' to 'const unsigned int', possible loss of data
C:\Users\sanchroy\Desktop\MTECH\deep learning\llama.cpp\tools\cvector-generator\cvector-generator.cpp(103): warning C4244: 'initializing': conversion from 'int64_t' to 'size_t', possible loss of data
C:\Users\sanchroy\Desktop\MTECH\deep learning\llama.cpp\tools\cvector-generator\cvector-generator.cpp(108): warning C4244: 'argument': conversion from 'float' to 'const unsigned int', possible loss of data
C:\Users\sanchroy\Desktop\MTECH\deep learning\llama.cpp\tools\cvector-generator\cvector-generator.cpp(119): warning C4244: 'initializing': conversion from 'int64_t' to 'int', possible loss of data
C:\Users\sanchroy\Desktop\MTECH\deep learning\llama.cpp\tools\cvector-generator\cvector-generator.cpp(129): warning C4305: 'argument': truncation from 'double' to 'float'
C:\Users\sanchroy\Desktop\MTECH\deep learning\llama.cpp\tools\cvector-generator\cvector-generator.cpp(137): warning C4244: 'initializing': conversion from 'int64_t' to 'int', possible loss of data
[ 99%] Linking CXX executable ..\..\bin\llama-cvector-generator.exe
[ 99%] Built target llama-cvector-generator
[ 99%] Building CXX object tools/export-lora/CMakeFiles/llama-export-lora.dir/export-lora.cpp.obj
export-lora.cpp
C:\Users\sanchroy\Desktop\MTECH\deep learning\llama.cpp\common\common.h(259): warning C4305: 'initializing': truncation from 'double' to 'float'
C:\Users\sanchroy\Desktop\MTECH\deep learning\llama.cpp\common\common.h(269): warning C4244: 'argument': conversion from 'const unsigned int' to 'float', possible loss of data
C:\Users\sanchroy\Desktop\MTECH\deep learning\llama.cpp\tools\export-lora\export-lora.cpp(22): warning C4244: 'initializing': conversion from 'int64_t' to 'int', possible loss of data
C:\Users\sanchroy\Desktop\MTECH\deep learning\llama.cpp\tools\export-lora\export-lora.cpp(27): warning C4244: 'initializing': conversion from 'int64_t' to 'int', possible loss of data
C:\Users\sanchroy\Desktop\MTECH\deep learning\llama.cpp\tools\export-lora\export-lora.cpp(319): warning C4244: 'argument': conversion from 'int64_t' to 'const unsigned int', possible loss of data
[100%] Linking CXX executable ..\..\bin\llama-export-lora.exe
[100%] Built target llama-export-lora

C:\Users\sanchroy\Desktop\MTECH\deep learning\llama.cpp>
```

```
C:\Users\sanchroy\Desktop\MTECH\deep learning\llama.cpp\build\bin>.\llama-bench.exe --help
warning: asserts enabled, performance may be affected
warning: debug build, performance may be affected
register_backend: registered backend CPU (1 devices)
register_device: registered device CPU (13th Gen Intel(R) Core(TM) i7-13800H)
load_backend: failed to find ggml_backend_init in C:\Users\sanchroy\Desktop\MTECH\deep learning\llama.cpp\build\bin\ggml-cpu.dll
usage: .\llama-bench.exe [options]

options:
  -h, --help
  --numa <distribute|isolate|numactl>      numa mode (default: disabled)
  -r, --repetitions <n>                      number of times to repeat each test (default: 5)
  --prio <-1|0|1|2|3>                         process/thread priority (default: 0)
  --delay <0...N> (seconds)                   delay between each test (default: 0)
  -o, --output <csv|json|jsonl|md|sql>        output format printed to stdout (default: md)
  -oe, --output-err <csv|json|jsonl|md|sql>    output format printed to stderr (default: none)
  --list-devices                                list available devices and exit
  -v, --verbose                                 verbose output
  --progress                                    print test progress indicators
  --no-warmup                                   skip warmup runs before benchmarking

test parameters:
  -m, --model <filename>                      (default: models/7B/ggml-model-q4_0.gguf)
  -p, --n-prompt <n>                           (default: 512)
  -n, --n-gen <n>                             (default: 128)
  -pg <pp,tgt>                                (default: )
  -d, --n-depth <n>                            (default: 0)
  -b, --batch-size <n>                          (default: 2048)
  -ub, --ubatch-size <n>                        (default: 512)
  -ctk, --cache-type-k <t>                     (default: f16)
  -ctv, --cache-type-v <t>                     (default: f16)
  -t, --threads <n>                            (default: 14)
  -C, --cpu-mask <hex,hex>                     (default: 0x0)
  --cpu-strict <0|1>                           (default: 0)
  --poll <0...100>                            (default: 50)
  -ngl, --n-gpu-layers <n>                     (default: 99)
  -ncmoe, --n-cpu-moe <n>                      (default: 0)
  -sm, --split-mode <none|layer|row>           (default: layer)
  -mg, --main-gpu <i>                           (default: 0)
  -nkvo, --no-kv-offload <0|1>                 (default: 0)
```

Task 2 – Setting up GPT-2 Medium Model

Steps to Execute:

1. **Download Model:** Cloning the GPT-2 Medium model repository directly from Hugging Face.

git clone <https://huggingface.co/openai-community/gpt2-medium>

2. Need to install other dependencies torch and transformer to make it work

3. **Convert to GGUF:** llama.cpp cannot load HuggingFace models directly.
It requires a special optimized format called GGUF (GGML Unified Format).

```
python3 convert_hf_to_gguf.py gpt2-medium --outfile gpt2-medium.gguf
```

4. **Run Sanity Benchmark:** A sanity benchmark on the converted GPT-2 Medium model using the llama-bench

```
-m ..\..\..\gpt2-medium.gguf : specifies the path to the converted model file.  
-p 0 : sets the prompt index to 0.  
-n 256 : runs the benchmark for 256 tokens to test model performance.
```

```
./llama-bench -m gpt2-medium.gguf -p 0 -n 256
```

Snapshots

```
C:\Users\sanchroy\Desktop\MTECH\deep learning>git clone https://huggingface.co/openai-community/gpt2-medium
Cloning into 'gpt2-medium'...
remote: Enumerating objects: 76, done.
remote: Total 76 (delta 0), reused 0 (delta 0), pack-reused 76 (from 1)
Unpacking objects: 100% (76/76), 1.65 MiB | 2.09 MiB/s, done.
Updating files: 100% (24/24), done.
```

```
C:\Users\sanchroy\Desktop\MTECH\deep learning\llama.cpp>pip install torch torchvision torchaudio
Defaulting to user installation because normal site-packages is not writeable
Requirement already satisfied: torch in c:\users\sanchroy\appdata\local\packages\pythonsoftwarefoundation.python.3.13_qbz5n2kfra8p0\localcache\local-packages\python313\site-packages (2.9.0)
Requirement already satisfied: torchvision in c:\users\sanchroy\appdata\local\packages\pythonsoftwarefoundation.python.3.13_qbz5n2kfra8p0\localcache\local-packages\python313\site-packages (0.24.0)
Collecting torchaudio
  Downloading torchaudio-2.9.0-cp313-cp313-win_amd64.whl.metadata (6.9 kB)
Requirement already satisfied: filelock in c:\users\sanchroy\appdata\local\packages\pythonsoftwarefoundation.python.3.13_qbz5n2kfra8p0\localcache\local-packages\python313\site-packages (from torch) (3.18.0)
Requirement already satisfied: typing-extensions>=4.10.0 in c:\users\sanchroy\appdata\local\packages\pythonsoftwarefoundation.python.3.13_qbz5n2kfra8p0\localcache\local-packages\python313\site-packages (from torch) (4.15.0)
Requirement already satisfied: sympy>=1.13.3 in c:\users\sanchroy\appdata\local\packages\pythonsoftwarefoundation.python.3.13_qbz5n2kfra8p0\localcache\local-packages\python313\site-packages (from torch) (1.14.0)
Requirement already satisfied: networkx>=2.5.1 in c:\users\sanchroy\appdata\local\packages\pythonsoftwarefoundation.python.3.13_qbz5n2kfra8p0\localcache\local-packages\python313\site-packages (from torch) (3.4.2)
Requirement already satisfied: jinja2 in c:\users\sanchroy\appdata\local\packages\pythonsoftwarefoundation.python.3.13_qbz5n2kfra8p0\localcache\local-packages\python313\site-packages (from torch) (3.1.6)
Requirement already satisfied: fsspec>=0.8.5 in c:\users\sanchroy\appdata\local\packages\pythonsoftwarefoundation.python.3.13_qbz5n2kfra8p0\localcache\local-packages\python313\site-packages (from torch) (2025.5.1)
Requirement already satisfied: setuptools in c:\users\sanchroy\appdata\local\packages\pythonsoftwarefoundation.python.3.13_qbz5n2kfra8p0\localcache\local-packages\python313\site-packages (from torch) (80.8.0)
Requirement already satisfied: numpy in c:\users\sanchroy\appdata\local\packages\pythonsoftwarefoundation.python.3.13_qbz5n2kfra8p0\localcache\local-packages\python313\site-packages (from torchvision) (2.2.4)
Requirement already satisfied: pillow!=8.3.* >=5.3.0 in c:\users\sanchroy\appdata\local\packages\pythonsoftwarefoundation.python.3.13_qbz5n2kfra8p0\localcache\local-packages\python313\site-packages (from torchvision) (11.1.0)
Requirement already satisfied: mpmath<1.4,>=1.1.0 in c:\users\sanchroy\appdata\local\packages\pythonsoftwarefoundation.python.3.13_qbz5n2kfra8p0\localcache\local-packages\python313\site-packages (from sympy>=1.13.3->torch) (1.3.0)
Requirement already satisfied: MarkupSafe>=2.0 in c:\users\sanchroy\appdata\local\packages\pythonsoftwarefoundation.python.3.13_qbz5n2kfra8p0\localcache\local-packages\python313\site-packages (from jinja2->torch) (3.0.2)
Downloading torchaudio-2.9.0-cp313-cp313-win_amd64.whl (665 kB)
665.3/665.3 kB 4.3 MB/s 0:00:00
Installing collected packages: torchaudio
Successfully installed torchaudio-2.9.0
[notice] A new release of pip is available: 25.2 -> 25.3
[notice] To update, run: C:\Users\sanchroy\AppData\Local\Microsoft\WindowsApps\PythonSoftwareFoundation.Python.3.13_qbz5n2kfra8p0\python.exe -m pip install --upgrade pip
```

```
C:\Users\sanchroy\Desktop\MTech\deep learning\llama.cpp>pip install transformers huggingface_hub tqdm sentencepiece
Defaulting to user installation because normal site-packages is not writeable
Collecting transformers
  Downloading transformers-4.57.1-py3-none-any.whl.metadata (43 kB)
Collecting huggingface_hub
  Downloading huggingface_hub-1.0.1-py3-none-any.whl.metadata (13 kB)
Requirement already satisfied: tqdm in c:\users\sanchroy\appdata\local\packages\pythonsoftwarefoundation.python.3.13_qbz5n2kfra8p0\localcache\local-packages\python313\site-packages (4.67.1)
Collecting sentencepiece
  Downloading sentencepiece-0.2.1-cp313-cp313-win_amd64.whl.metadata (10 kB)
Requirement already satisfied: numpy>=1.17 in c:\users\sanchroy\appdata\local\packages\pythonsoftwarefoundation.python.3.13_qbz5n2kfra8p0\localcache\local-packages\python313\site-packages (from transformers) (3.18.0)
Collecting huggingface_hub
  Downloading huggingface_hub-0.36.0-py3-none-any.whl.metadata (14 kB)
Requirement already satisfied: packaging>=20.0 in c:\users\sanchroy\appdata\local\packages\pythonsoftwarefoundation.python.3.13_qbz5n2kfra8p0\localcache\local-packages\python313\site-packages (from transformers) (2.2.4)
Requirement already satisfied: pyyaml>=5.1 in c:\users\sanchroy\appdata\local\packages\pythonsoftwarefoundation.python.3.13_qbz5n2kfra8p0\localcache\local-packages\python313\site-packages (from transformers) (6.0.3)
Collecting regex<=2019.12.17 (from transformers)
  Downloading regex-2025.10.23-cp313-cp313-win_amd64.whl.metadata (41 kB)
Requirement already satisfied: requests in c:\users\sanchroy\appdata\local\packages\pythonsoftwarefoundation.python.3.13_qbz5n2kfra8p0\localcache\local-packages\python313\site-packages (from transformers) (2.32.5)
Collecting tokenizers<=0.23.0 >=0.22.0 (from transformers)
  Downloading tokenizers-0.22.1-cp39-ab13-win_amd64.whl.metadata (6.9 kB)
Collecting safetensors>=0.4.3 (from transformers)
  Downloading safetensors-0.6.2-cp38-ab13-win_amd64.whl.metadata (4.1 kB)
Requirement already satisfied: fsspec>=2023.5.0 in c:\users\sanchroy\appdata\local\packages\pythonsoftwarefoundation.python.3.13_qbz5n2kfra8p0\localcache\local-packages\python313\site-packages (from huggingface_hub) (2025.5.1)
Requirement already satisfied: typing_extensions>=3.7.4.3 in c:\users\sanchroy\appdata\local\packages\pythonsoftwarefoundation.python.3.13_qbz5n2kfra8p0\localcache\local-packages\python313\site-packages (from huggingface_hub) (4.15.0)
Requirement already satisfied: colorama in c:\users\sanchroy\appdata\local\packages\pythonsoftwarefoundation.python.3.13_qbz5n2kfra8p0\localcache\local-packages\python313\site-packages (from tqdm) (0.4.6)
Requirement already satisfied: charset_normalizer<4,>=2 in c:\users\sanchroy\appdata\local\packages\pythonsoftwarefoundation.python.3.13_qbz5n2kfra8p0\localcache\local-packages\python313\site-packages (from requests->transformers) (3.4.4)
Requirement already satisfied: idna<4,>=2.5 in c:\users\sanchroy\appdata\local\packages\pythonsoftwarefoundation.python.3.13_qbz5n2kfra8p0\localcache\local-
```

```

  Downloading safetensors-0.6.2-cp38-win_amd64.whl.metadata (4.1 kB)
Requirement already satisfied: fsspec>=2023.5.0 in c:\users\sanchroy\appdata\local\packages\pythonsoftwarefoundation.python.3.13_qbz5n2kfra8p0\localcache\local-packages\python313\site-packages (from huggingface_hub) (2025.5.1)
Requirement already satisfied: typing-extensions>=3.7.4.3 in c:\users\sanchroy\appdata\local\packages\pythonsoftwarefoundation.python.3.13_qbz5n2kfra8p0\localcache\local-packages\python313\site-packages (from huggingface_hub) (4.15.0)
Requirement already satisfied: colorama in c:\users\sanchroy\appdata\local\packages\pythonsoftwarefoundation.python.3.13_qbz5n2kfra8p0\localcache\local-packages\python313\site-packages (from tqdm) (0.4.6)
Requirement already satisfied: charset_normalizer<4,>=2 in c:\users\sanchroy\appdata\local\packages\pythonsoftwarefoundation.python.3.13_qbz5n2kfra8p0\localcache\local-packages\python313\site-packages (from requests->transformers) (3.4.4)
Requirement already satisfied: idna<4,>=2.5 in c:\users\sanchroy\appdata\local\packages\pythonsoftwarefoundation.python.3.13_qbz5n2kfra8p0\localcache\local-packages\python313\site-packages (from requests->transformers) (3.11)
Requirement already satisfied: urllib3<3,>=1.21.1 in c:\users\sanchroy\appdata\local\packages\pythonsoftwarefoundation.python.3.13_qbz5n2kfra8p0\localcache\local-packages\python313\site-packages (from requests->transformers) (2.5.0)
Requirement already satisfied: certifi>=2017.4.17 in c:\users\sanchroy\appdata\local\packages\pythonsoftwarefoundation.python.3.13_qbz5n2kfra8p0\localcache\local-packages\python313\site-packages (from requests->transformers) (2025.10.5)
  Downloading transformers-4.57.1-py3-none-any.whl (12.0 MB)
  └── transformers-4.57.1-py3-none-any.whl (12.0 MB) 12.0/12.0 MB 5.0 MB/s 0:00:02
  Downloading huggingface_hub-0.36.0-py3-none-any.whl (566 kB)
  └── huggingface_hub-0.36.0-py3-none-any.whl (566 kB) 566.1/566.1 kB 4.4 MB/s 0:00:00
  Downloading tokenizers-0.22.1-cp39-abi3-win_amd64.whl (2.7 MB)
  └── tokenizers-0.22.1-cp39-abi3-win_amd64.whl (2.7 MB) 2.7/2.7 MB 5.2 MB/s 0:00:00
  Downloading sentencepiece-0.2.1-cp313-cp313-win_amd64.whl (1.1 MB)
  └── sentencepiece-0.2.1-cp313-cp313-win_amd64.whl (1.1 MB) 1.1/1.1 MB 4.6 MB/s 0:00:00
  Downloading regex-2025.10.23-cp313-cp313-win_amd64.whl (276 kB)
  Downloading safetensors-0.6.2-cp38-abi3-win_amd64.whl (320 kB)
  Downloading collected packages: sentencepiece, safetensors, regex, huggingface_hub, tokenizers, transformers
  └── sentencepiece-0.2.1-cp313-cp313-win_amd64.whl [huggingface_hub] WARNING: The scripts hf.exe, huggingface-cl.exe and tiny-agents.exe are installed in 'C:\Users\sanchroy\AppData\Local\ Packages\PythonSoftwareFoundation.Python.3.13_qbz5n2kfra8p0\LocalCache\local-packages\Python313\Scripts' which is not on PATH.
  Consider adding this directory to PATH or, if you prefer to suppress this warning, use --no-warn-script-location.
  └── transformers-4.57.1-py3-none-any.whl [transformers] WARNING: The scripts transformers-cl.exe and transformers.exe are installed in 'C:\Users\sanchroy\AppData\Local\ Packages\PythonSoftwareFoundation.Python.3.13_qbz5n2kfra8p0\LocalCache\local-packages\Python313\Scripts' which is not on PATH.
  Consider adding this directory to PATH or, if you prefer to suppress this warning, use --no-warn-script-location.
Successfully installed huggingface_hub-0.36.0 regex-2025.10.23 safetensors-0.6.2 sentencepiece-0.2.1 tokenizers-0.22.1 transformers-4.57.1

[notice] A new release of pip is available: 25.2 → 25.3
[notice] To update, run: C:\Users\sanchroy\AppData\Local\Microsoft\WindowsApps\PythonSoftwareFoundation.Python.3.13_qbz5n2kfra8p0\python.exe -m pip install --upgrade pip
C:\Users\sanchroy\Desktop\MTech\deep learning\llama.cpp>
C:\Users\sanchroy\Desktop\MTech\deep learning\llama.cpp>
C:\Users\sanchroy\Desktop\MTech\deep learning\llama.cpp>

C:\Users\sanchroy\Desktop\MTech\deep learning\llama.cpp>python3 convert_hf_to_gguf.py "C:\Users\sanchroy\Desktop\MTech\deep learning\gpt2-medium" --outfile gpt2-medium.gguf
INFO:hf-to-gguf:Loading model: gpt2-medium
INFO:hf-to-gguf:Model architecture: GPT2LMHeadModel
INFO:hf-to-gguf:gguf: indexing model part 'model.safetensors'
INFO:gguf:gguf_writer:gguf: This GGUF file is for Little Endian only
INFO:hf-to-gguf:Exporting model...
INFO:hf-to-gguf:blk.0.attn_qkv.bias, torch.float32 --> F32, shape = {3072}
INFO:hf-to-gguf:blk.0.attn_qkv.weight, torch.float32 --> F16, shape = {1024, 3072}
INFO:hf-to-gguf:blk.0.attn_output.bias, torch.float32 --> F32, shape = {1024}
INFO:hf-to-gguf:blk.0.attn_output.weight, torch.float32 --> F16, shape = {1024, 1024}
INFO:hf-to-gguf:blk.0.attn_norm.bias, torch.float32 --> F32, shape = {1024}
INFO:hf-to-gguf:blk.0.attn_norm.weight, torch.float32 --> F32, shape = {1024}
INFO:hf-to-gguf:blk.0.ffn_norm.bias, torch.float32 --> F32, shape = {1024}
INFO:hf-to-gguf:blk.0.ffn_norm.weight, torch.float32 --> F32, shape = {1024}
INFO:hf-to-gguf:blk.0.ffn_up.bias, torch.float32 --> F32, shape = {4096}
INFO:hf-to-gguf:blk.0.ffn_up.weight, torch.float32 --> F16, shape = {1024, 4096}
INFO:hf-to-gguf:blk.0.ffn_down.bias, torch.float32 --> F32, shape = {1024}
INFO:hf-to-gguf:blk.0.ffn_down.weight, torch.float32 --> F16, shape = {4096, 1024}
INFO:hf-to-gguf:blk.1.attn_qkv.bias, torch.float32 --> F32, shape = {3072}
INFO:hf-to-gguf:blk.1.attn_qkv.weight, torch.float32 --> F16, shape = {1024, 3072}
INFO:hf-to-gguf:blk.1.attn_output.bias, torch.float32 --> F32, shape = {1024}
INFO:hf-to-gguf:blk.1.attn_output.weight, torch.float32 --> F16, shape = {1024, 1024}
INFO:hf-to-gguf:blk.1.attn_norm.bias, torch.float32 --> F32, shape = {1024}
INFO:hf-to-gguf:blk.1.attn_norm.weight, torch.float32 --> F32, shape = {1024}
INFO:hf-to-gguf:blk.1.ffn_norm.bias, torch.float32 --> F32, shape = {1024}
INFO:hf-to-gguf:blk.1.ffn_norm.weight, torch.float32 --> F32, shape = {1024}
INFO:hf-to-gguf:blk.1.ffn_up.bias, torch.float32 --> F32, shape = {4096}
INFO:hf-to-gguf:blk.1.ffn_up.weight, torch.float32 --> F16, shape = {1024, 4096}
INFO:hf-to-gguf:blk.1.ffn_down.bias, torch.float32 --> F32, shape = {1024}
INFO:hf-to-gguf:blk.1.ffn_down.weight, torch.float32 --> F16, shape = {4096, 1024}
INFO:hf-to-gguf:blk.10.attn_qkv.bias, torch.float32 --> F32, shape = {3072}
INFO:hf-to-gguf:blk.10.attn_qkv.weight, torch.float32 --> F16, shape = {1024, 3072}
INFO:hf-to-gguf:blk.10.attn_output.bias, torch.float32 --> F32, shape = {1024}
INFO:hf-to-gguf:blk.10.attn_output.weight, torch.float32 --> F16, shape = {1024, 1024}
INFO:hf-to-gguf:blk.10.attn_norm.bias, torch.float32 --> F32, shape = {1024}
INFO:hf-to-gguf:blk.10.attn_norm.weight, torch.float32 --> F32, shape = {1024}

```

```

INFO:hf-to-gguf:blk.8.ffn_up.bias, torch.float32 --> F32, shape = {4096}
INFO:hf-to-gguf:blk.8.ffn_up.weight, torch.float32 --> F16, shape = {1024, 4096}
INFO:hf-to-gguf:blk.8.ffn_down.bias, torch.float32 --> F32, shape = {1024}
INFO:hf-to-gguf:blk.8.ffn_down.weight, torch.float32 --> F16, shape = {4096, 1024}
INFO:hf-to-gguf:blk.9.attn_qkv.bias, torch.float32 --> F32, shape = {3072}
INFO:hf-to-gguf:blk.9.attn_qkv.weight, torch.float32 --> F16, shape = {1024, 3072}
INFO:hf-to-gguf:blk.9.attn_output.bias, torch.float32 --> F32, shape = {1024}
INFO:hf-to-gguf:blk.9.attn_output.weight, torch.float32 --> F16, shape = {1024, 1024}
INFO:hf-to-gguf:blk.9.attn_norm.bias, torch.float32 --> F32, shape = {1024}
INFO:hf-to-gguf:blk.9.ffn_norm.bias, torch.float32 --> F32, shape = {1024}
INFO:hf-to-gguf:blk.9.ffn_norm.weight, torch.float32 --> F32, shape = {1024}
INFO:hf-to-gguf:blk.9.ffn_up.bias, torch.float32 --> F32, shape = {4096}
INFO:hf-to-gguf:blk.9.ffn_up.weight, torch.float32 --> F16, shape = {1024, 4096}
INFO:hf-to-gguf:blk.9.ffn_down.bias, torch.float32 --> F32, shape = {1024}
INFO:hf-to-gguf:blk.9.ffn_down.weight, torch.float32 --> F16, shape = {4096, 1024}
INFO:hf-to-gguf:output_norm.bias, torch.float32 --> F32, shape = {1024}
INFO:hf-to-gguf:output_norm.weight, torch.float32 --> F32, shape = {1024}
INFO:hf-to-gguf:position_embd.weight, torch.float32 --> F32, shape = {1024, 1024}
INFO:hf-to-gguf:token_embd.weight, torch.float32 --> F16, shape = {1024, 50257}
INFO:hf-to-gguf:Set meta model
INFO:hf-to-gguf:Set model parameters
INFO:hf-to-gguf:Set model quantization version
INFO:hf-to-gguf:Set model tokenizer
INFO:gguf.vocab:Adding 50000 merge(s).
INFO:gguf.vocab:Setting special token type bos to 50256
INFO:gguf.vocab:Setting special token type eos to 50256
INFO:gguf.vocab:Setting special token type unk to 50256
INFO:gguf.gguf_writer:Writing the following files:
INFO:gguf.gguf_writer:gpt2-medium.gguf: n_tensors = 292, total_size = 712.4M
Writing: 100% [██████████] 712M/712M [00:15<00:00, 45.8Mbyte/s]
INFO:hf-to-gguf:Model successfully exported to gpt2-medium.gguf

```

```

::\Users\sanchroy\Desktop\MTECH\deep learning\llama.cpp\build\bin\Release>llama-bench.exe -m "C:\Users\sanchroy\Desktop\MTECH\deep learning\llama.cpp\gpt2-medium.gguf" -p 0 -n 256
| model           | size | params | backend | threads | test | t/s |
|-----:|-----:|-----:|-----:|-----:|-----:|-----:|
| gpt2 0.4B F16 | 679.38 MiB | 354.82 M | CPU     | 14      | tg256 | 55.19 ± 1.27 |
build: 1ae74882f (6913)
::\Users\sanchroy\Desktop\MTECH\deep learning\llama.cpp\build\bin\Release>

```

Explanation:

Field	Value	Explanation
Model	gpt2 0.4B F16	GPT-2 Medium (354M parameters) loaded in FP16 GGUF format. “0.4B” refers to ~0.354B parameters.
Model Size	679.38 MiB	Size of the FP16 GGUF file. FP16 doubles the memory compared to quantized models but gives higher accuracy.
Parameters	354.82 M	Number of trainable weights in GPT-2 Medium. llama.cpp reads this directly from GGUF metadata.
Backend	CPU	The benchmark used pure CPU execution (no GPU acceleration).
Threads	14	llama.cpp used 14 CPU threads for parallel inference. This is close to the physical + logical CPU cores on my system.
Test	tg256	Built-in benchmark “token generation 256” which measures generation speed for 256 tokens.

Field	Value	Explanation
Throughput (t/s)	55.19 ± 1.27	The most important metric. CPU generates ~55 tokens/sec with 14 threads—strong performance for GPT-2 Medium in FP16 mode.

Task 3 – Naive Execution (No Parallelism)

- A. Re-build llama.cpp by disabling all acceleration flags.
- B. Run Benchmark: After compilation, run the benchmark in single-thread *mode*:

Deliverable:

Submit benchmark results from the single-thread run.

Solution:

- A. Rebuild llama.cpp with all acceleration flags disabled

Steps to Execute:

```
cmake -B build -DGGML_CPU_GENERIC=ON -DGGML_NATIVE=OFF -DGGML_AVX=OFF -  
DGGML_AVX2=OFF -DGGML_AVX512=OFF -DGGML_SSE42=OFF -DGGML_F16C=OFF -  
DGGML_FMA=OFF
```

-B build: Generates build files in the *build* directory.

-DGGML_CPU_GENERIC=ON: Enables a generic CPU build for all processors.

-DGGML_NATIVE=OFF: Disables host-specific optimizations.

-DGGML_AVX/AVX2/AVX512=OFF: Turns off AVX instruction sets for parallel processing.

-DGGML_SSE42=OFF: Disables SSE4.2 optimizations.

-DGGML_F16C=OFF: Disables half-precision (16-bit) float support.

-DGGML_FMA=OFF: Disables Fused Multiply-Add arithmetic optimization.

Details:

This configuration forces the program to compile without any hardware acceleration, ensuring that it runs purely on basic CPU instructions. It is useful for benchmarking and testing the performance of the program under non-optimized conditions.

Run the benchmark in single-thread mode after compilation.

-m ..\..\gpt2-medium.gguf : specifies the path to the converted model file.

-p 0 : sets the prompt index to 0 .

-n 256 : runs the benchmark for 256 tokens for model performance.

-t 1: Runs the benchmark on a single CPU thread

Snapshots

```
sanchari@ubuntults:~/llama.cpp$ sanchari@ubuntults:~/llama.cpp$ cmake -B build -DGGML_CPU_GENERIC=ON -DGGML_NATIVE=OFF -DGGML_AVX=OFF -DGGML_AVX2=OFF -DGGML_AVX512=OFF -DGGML_SSE42=OFF -DGGML_F16C=OFF -DGGML_FMA=OFF CMAKE_BUILD_TYPE=Release
-- ccache found, compilation results will be cached. Disable with GGML_CCACHE=OFF.
-- CMAKE SYSTEM PROCESSOR: x86_64
-- GGML SYSTEM ARCH: x86
-- Including CPU backend
-- x86 detected
-- Adding CPU backend variant ggml-cpu:
-- ggml version: 0.9.4
-- ggml commit: 9b17d7d4ab
-- Configuration done (0.5s)
-- Generating done (0.5s)
-- Build files have been written to: /home/sanchari/llama.cpp/build
sanchari@ubuntults:~/llama.cpp$ cmake --build build --config Release
[ 3%] Built target ggml-base
[ 7%] Built target ggml-cpu
[ 7%] Built target ggml
[ 43%] Built target llama
[ 43%] Built target build_info
[ 47%] Built target common
[ 48%] Built target cpp-httplib
[ 49%] Built target test-tokenizer-0
[ 50%] Built target test-sampling
[ 51%] Built target test-grammar-parser
[ 52%] Built target test-grammar-integration
[ 53%] Built target test-llama-grammar
[ 54%] Built target test-chat
[ 54%] Built target test-json-schema-to-grammar
[ 55%] Built target test-quantize-stats
[ 55%] Built target test-gbnf-validator
[ 56%] Built target test-tokenizer-l-bpe
[ 56%] Built target test-tokenizer-l-spm
[ 57%] Built target test-chat-parser
[ 58%] Built target test-chat-template
[ 59%] Built target test-json-partial
[ 60%] Built target test-log
[ 61%] Built target test-regex-partial
[ 62%] Built target test-thread-safety
[ 62%] Built target test-arg-parser
[ 63%] Built target test-quant
[ 64%] Built target test-gnuf
[ 65%] Built target test-backend-ops
[ 66%] Built target test-model-load-cancel
[ 67%] Built target test-autorelease
[ 67%] Built target test-barrier
[ 68%] Built target test-quantize-fns
[ 68%] Built target test-quantize-n erf
```

```

[ 84%] Built target llama-speculative-simple
[ 84%] Built target llama-gen-docs
[ 85%] Built target llama-finetune
[ 86%] Built target llama-diffusion-cli
[ 87%] Built target llama-logits
[ 88%] Built target llama-convert-llama2c-to-ggml
[ 88%] Built target llama-vdot
[ 89%] Built target llama-qdot
[ 90%] Built target llama-batched-bench
[ 91%] Built target llama-gguf-split
[ 91%] Built target llama-imatrix
[ 91%] Built target llama-bench
[ 92%] Built target llama-cli
[ 92%] Built target llama-perplexity
[ 92%] Built target llama-quantize
[ 93%] Built target llama-server
[ 93%] Built target llama-run
[ 94%] Built target llama-tokenize
[ 95%] Built target llama-tts
[ 96%] Built target llama-llava-cli
[ 97%] Built target llama-gemma3-cli
[ 98%] Built target llama-minicpmv-cli
[ 99%] Built target llama-qwen2vl-cli
[100%] Built target llama-mtmd-cli
[100%] Built target llama-cvector-generator
[100%] Built target llama-export-lora
sanchari@ubuntults:~/llama.cpp$ ./build/bin/llama-bench -m ./gpt2-medium.gguf -p 0 -n 256 -t 1
./build/bin/llama-bench: error while loading shared libraries: libsvm1.so: cannot open shared object file: No such file or directory
sanchari@ubuntults:~/llama.cpp$ source /opt/intel/oneapi/setvars.sh

:: initializing oneAPI environment ...
-bash: BASH_VERSION = 5.2.21(1)-release
args: Using ${@} for setvars.sh arguments:
:: compiler -- latest
:: debugger -- latest
:: dev-utilities -- latest
:: dpl -- latest
:: mkl -- latest
:: tbb -- latest
:: umf -- latest
:: oneAPI environment initialized ::

sanchari@ubuntults:~/llama.cpp$ ./build/bin/llama-bench -m ./gpt2-medium.gguf -p 0 -n 256 -t 1
+-----+-----+-----+-----+-----+-----+-----+
| model | size | params | backend | threads | test | t/s |
+-----+-----+-----+-----+-----+-----+-----+
| gpt2 0.4B F16 | 679.38 MiB | 354.82 M | CPU | 1 | tg256 | 1.40 ± 0.01 |

build: 9b17d74ab (7062)
sanchari@ubuntults:~/llama.cpp$ █

```

[Explanation](#)

Field	Value	Explanation
Model	gpt2 0.4B F16	GPT-2 Medium in FP16 GGUF format (\approx 354M parameters).
Model Size	679.38 MiB	FP16 GGUF file loaded by llama.cpp.
Parameters	354.82 M	Total number of FP16 weights in the model.
Backend	CPU (Naive Build)	No SIMD, no MKL, no OpenBLAS
Threads	1	Task 3 requires <i>single-thread execution</i> to measure true naieve performance.
Test	tg256	Standard 256-token generation benchmark used by llama.cpp.
Throughput (t/s)	1.40 ± 0.01 tokens/sec	Extremely low speed due to no parallelism or CPU vectorization.

Task 4 – Default Execution

- A. Build with default settings
- B. Run Benchmark

Deliverables:

Submit a screenshot of the benchmark output for the single-threaded execution.

[Solution:](#)

Steps to Execute:

cmake -B build

-B build: Generates build files in the *build* directory.

Details:

Building with default settings.

Run the benchmark in single-thread mode after compilation.

-m ..\..\gpt2-medium.gguf : specifies the path to the converted model file.

-p 0 : sets the prompt index to 0 .

-n 256 : runs the benchmark for 256 tokens for model performance.

-t 1:

```
sanchari@ubuntults:~/llama.cpp$ cmake -B build
CMAKE_BUILD_TYPE=Release
-- ccache found, compilation results will be cached. Disable with GGML_CCACHE=OFF.
-- CMAKE_SYSTEM_PROCESSOR: x86_64
-- GGML_SYSTEM_ARCH: x86
-- Including CPU backend
-- x86 detected
-- Adding CPU backend variant ggml-cpu:
-- ggml version: 0.9.4
-- ggml commit: 9b17d74ab
-- Configuring done (0.6s)
-- Generating done (0.5s)
-- Build files have been written to: /home/sanchari/llama.cpp/build
sanchari@ubuntults:~/llama.cpp$ cmake --build build --config Release
[ 3%] Built target ggml-base
[ 7%] Built target ggml-cpu
[ 7%] Built target ggml
[ 43%] Built target llama
[ 43%] Built target build_info
[ 47%] Built target common
[ 48%] Built target cpp-httplib
[ 49%] Built target test-tokenizer-0
[ 50%] Built target test-sampling
[ 51%] Built target test-grammar-parser
[ 52%] Built target test-grammar-integration
[ 53%] Built target test-llama-grammar
[ 54%] Built target test-chat
[ 54%] Built target test-json-schema-to-grammar
[ 55%] Built target test-quantize-stats
[ 55%] Built target test-qbnf-validator
[ 56%] Built target test-tokenizer-l-bpe
[ 56%] Built target test-tokenizer-l-spm
[ 57%] Built target test-chat-parser
[ 58%] Built target test-chat-template
[ 59%] Built target test-json-partial
[ 60%] Built target test-log
[ 61%] Built target test-regex-partial
[ 61%] Built target test-thread-safety
[ 62%] Built target test-arg-parser
[ 63%] Built target test-opt
[ 64%] Built target test-gguf
[ 65%] Built target test-backend-ops
[ 66%] Built target test-model-load-cancel
[ 67%] Built target test-autorelease
[ 67%] Built target test-barrier
[ 68%] Built target test-quantize-fns
[ 68%] Built target test-quantize-perf
[ 69%] Built target test-rope
[ 71%] Built target mtmd
```

```

[ 76%] Built target xxhash
[ 76%] Built target sha1
[ 76%] Built target llama-gguf-hash
[ 77%] Built target llama-gguf
[ 77%] Built target llama-lookahead
[ 78%] Built target llama-lookup
[ 78%] Built target llama-lookup-create
[ 79%] Built target llama-lookup-merge
[ 79%] Built target llama-lookup-stats
[ 79%] Built target llama-parallel
[ 80%] Built target llama-passkey
[ 81%] Built target llama-retrieval
[ 82%] Built target llama-save-load-state
[ 83%] Built target llama-simple
[ 83%] Built target llama-simple-chat
[ 84%] Built target llama-speculative
[ 84%] Built target llama-speculative-simple
[ 84%] Built target llama-gen-docs
[ 85%] Built target llama-finetune
[ 86%] Built target llama-difusion-cli
[ 87%] Built target llama-logits
[ 88%] Built target llama-convert-llama2c-to-ggml
[ 88%] Built target llama-vdot
[ 89%] Built target llama-qdot
[ 90%] Built target llama-batched-bench
[ 91%] Built target llama-gguf-split
[ 91%] Built target llama-imatrix
[ 91%] Built target llama-bench
[ 92%] Built target llama-cli
[ 92%] Built target llama-perplexity
[ 92%] Built target llama-quantize
[ 93%] Built target llama-server
[ 93%] Built target llama-run
[ 94%] Built target llama-tokenize
[ 95%] Built target llama-tts
[ 96%] Built target llama-llava-cli
[ 97%] Built target llama-gemma3-cli
[ 98%] Built target llama-minicpmv-cli
[ 99%] Built target llama-qwenzvl-cli
[100%] Built target llama-mtmd-cli
[100%] Built target llama-cvector-generator
[100%] Built target llama-export-lora
sanchari@ubuntults:~/llama.cpp$ ./build/bin/llama-bench -m ./gpt2-medium.gguf -p 0 -n 256 -t 1
| model           | size       | params    | backend   | threads | test      | t/s |
| -----           | -----:     | -----:    | -----:   | -----:  | -----:    | -----: |
| gpt2 0.4B F16  | 679.38 MiB | 354.82 M | CPU       | 1        | tg256     | 1.40 ± 0.00 |

build: 9b17d74ab (7062)
sanchari@ubuntults:~/llama.cpp$ █

```

Explanation

Field	Value	Explanation
Model	gpt2 0.4B F16	GPT-2 Medium model in FP16 GGUF format (~354M parameters).
Model Size	679.38 MiB	Size of the loaded FP16 model file.
Parameters	354.82 M	Number of neural network weights in GPT-2 Medium.
Backend	CPU	Naive llama.cpp build—no BLAS, no MKL, no SIMD acceleration.
Threads	1	Task 4 is evaluated single-threaded to measure baseline performance.
Test	tg256	Benchmarks token-generation performance for 256 tokens.
Throughput (t/s)	1.40 ± 0.00 tokens/sec	Slow performance; reflects absence of SIMD or BLAS optimizations.

Task 5 – Near-Optimal Execution with Intel MKL

A. Rebuild with Intel MKL.

B. Run Benchmark Once built, run the benchmark with a single thread

Deliverables:

Submit a screenshot of the benchmark output for the single-threaded execution.

Solution:

Re-Building with Intel MKL settings

Commands:

-B build:

-DGGML_BLAS=ON
-DGGML_BLAS_VENDOR=Intel10_64lp
-DCMAKE_C_COMPILER=icx
-DCMAKE_CXX_COMPILER=icpx
-DGGML_NATIVE=ON

Generates build files in the build directory

Enables BLAS acceleration in GGML/llama.cpp.

Selects Intel MKL (LP64 interface).

Uses Intel oneAPI icx as the C compiler.

Uses Intel oneAPI icpx as the C++ compiler.

Enables CPU-specific optimizations (-march=native).

Run the benchmark in single-thread mode after compilation.

Screenshots

```
sanchari@ubuntults:~/llama.cpp$ cmake -B build -DGGML_BLAS=ON -DGGML_BLAS_VENDOR=Intel10_64lp \
-DCMAKE_C_COMPILER=icx -DCMAKE_CXX_COMPILER=icpx -DGGML_NATIVE=ON
CMAKE_BUILD_TYPE=Release
-- ccache found, compilation results will be cached. Disable with GGML_CCACHE=OFF.
-- CMAKE_SYSTEM_PROCESSOR: x86_64
-- GGML_SYSTEM_ARCH: x86
-- Including CPU backend
-- x86 detected
-- Adding CPU backend variant ggml-cpu: -march=native
-- Looking for sgemm_
-- Looking for sgemm_ - found
-- Found BLAS: /opt/intel/oneapi/mkl/2025.3/lib/libmkl_intel_lp64.so;/opt/intel/oneapi/mkl/2025.3/lib/libmkl_intel_thread.so;/opt/intel/oneapi/mkl/2025.3/lib/libmkl_i
/opt/intel/oneapi/compiler/2025.3/lib/libimomp5.so;-lm;-ldl
-- BLAS found, Libraries: /opt/intel/oneapi/mkl/2025.3/lib/libmkl_intel_lp64.so;/opt/intel/oneapi/mkl/2025.3/lib/libmkl_intel_thread.so;/opt/intel/oneapi/mkl/2025.
intel/oneapi/compiler/2025.3/lib/libimomp5.so;-lm;-ldl
-- Found PkgConfig: /usr/bin/pkg-config (found version "1.8.1")
-- Checking for module 'mkl-sdl'
--   Found mkl-sdl, version 2025.3
-- BLAS found, Includes: /opt/intel/oneapi/mkl/2025.3/lib/pkgconfig/.../include
-- Including BLAS backend
-- ggml version: 0.9.4
-- ggml commit: 9b17d74ab
-- Configuring done (1.5s)
-- Generating done (0.5s)
-- Build files have been written to: /home/sanchari/llama.cpp/build
sanchari@ubuntults:~/llama.cpp$ cmake --build build --config Release
[ 3%] Built target ggml-base
[ 3%] Building CXX object ggml/src/ggml-blas/CMakeFiles/ggml-blas.dir/ggml-blas.cpp.o
[ 3%] Linking CXX shared library ../../bin/libggml-blas.so
[ 3%] Built target ggml-blas
[ 4%] Building C object ggml/src/CMakeFiles/ggml-cpu.dir/ggml-cpu/ggml-cpu.c.o
[ 4%] Building CXX object ggml/src/CMakeFiles/ggml-cpu.dir/ggml-cpu/ggml-cpu.cpp.o
[ 4%] Building CXX object ggml/src/CMakeFiles/ggml-cpu.dir/ggml-cpu/repack.cpp.o
[ 4%] Building CXX object ggml/src/CMakeFiles/ggml-cpu.dir/ggml-cpu/hbm.cpp.o
[ 5%] Building C object ggml/src/CMakeFiles/ggml-cpu.dir/ggml-cpu/quants.c.o
[ 5%] Building CXX object ggml/src/CMakeFiles/ggml-cpu.dir/ggml-cpu/traits.cpp.o
[ 5%] Building CXX object ggml/src/CMakeFiles/ggml-cpu.dir/ggml-cpu/amx/amx.cpp.o
[ 6%] Building CXX object ggml/src/CMakeFiles/ggml-cpu.dir/ggml-cpu/amx/mmq.cpp.o
[ 6%] Building CXX object ggml/src/CMakeFiles/ggml-cpu.dir/ggml-cpu/binary-ops.cpp.o
[ 6%] Building CXX object ggml/src/CMakeFiles/ggml-cpu.dir/ggml-cpu/unary-ops.cpp.o
[ 6%] Building CXX object ggml/src/CMakeFiles/ggml-cpu.dir/ggml-cpu/vec.cpp.o
[ 7%] Building CXX object ggml/src/CMakeFiles/ggml-cpu.dir/ggml-cpu/ops.cpp.o
[ 7%] Building CXX object ggml/src/CMakeFiles/ggml-cpu.dir/ggml-cpu/llamofile/sgemm.cpp.o
[ 7%] Building C object ggml/src/CMakeFiles/ggml-cpu.dir/ggml-cpu/arch/x86/quants.c.o
[ 8%] Building CXX object ggml/src/CMakeFiles/ggml-cpu.dir/ggml-cpu/arch/x86/repack.cpp.o
[ 8%] Linking CXX shared library ../../bin/libggml-cpu.so
[ 8%] Built target ggml-cpu
[ 8%] Building CXX object ggml/src/CMakeFiles/ggml.dir/ggml-backend-reg.cpp.o
```

```

[ 91%] Built target llama-gguf-split
[ 91%] Building CXX object tools/imatrix/CMakeFiles/llama-imatrix.dir/imatrix.cpp.o
[ 92%] Linking CXX executable ../../bin/llama-imatrix
[ 92%] Built target llama-imatrix
[ 92%] Building CXX object tools/llama-bench/CMakeFiles/llama-bench.dir/llama-bench.cpp.o
[ 93%] Linking CXX executable ../../bin/llama-bench
[ 93%] Built target llama-bench
[ 93%] Building CXX object tools/main/CMakeFiles/llama-cli.dir/main.cpp.o
[ 93%] Linking CXX executable ../../bin/llama-cli
[ 93%] Built target llama-cli
[ 93%] Building CXX object tools/perplexity/CMakeFiles/llama-perplexity.dir/perplexity.cpp.o
[ 93%] Linking CXX executable ../../bin/llama-perplexity
[ 93%] Built target llama-perplexity
[ 93%] Building CXX object tools/quantize/CMakeFiles/llama-quantize.dir/quantize.cpp.o
[ 94%] Linking CXX executable ../../bin/llama-quantize
[ 94%] Built target llama-quantize
[ 94%] Building CXX object tools/server/CMakeFiles/llama-server.dir/server.cpp.o
[ 94%] Linking CXX executable ../../bin/llama-server
[ 95%] Built target llama-server
[ 95%] Building CXX object tools/run/CMakeFiles/llama-run.dir/run.cpp.o
[ 95%] Building CXX object tools/run/CMakeFiles/llama-run.dir/linenoise.cpp/linenoise.cpp.o
[ 96%] Linking CXX executable ../../bin/llama-run
[ 96%] Built target llama-run
[ 96%] Building CXX object tools/tokenize/CMakeFiles/llama-tokenize.dir/tokenize.cpp.o
[ 96%] Linking CXX executable ../../bin/llama-tokenize
[ 96%] Built target llama-tokenize
[ 96%] Building CXX object tools/tts/CMakeFiles/llama-tts.dir/tts.cpp.o
[ 97%] Linking CXX executable ../../bin/llama-tts
[ 97%] Built target llama-tts
[ 97%] Built target llama-llava-cli
[ 98%] Built target llama-gemm3-cli
[ 98%] Built target llama-minicpmv-cli
[ 98%] Built target llama-qwen2vl-cli
[ 98%] Building CXX object tools/mtmd/CMakeFiles/llama-mtmd-cli.dir/mtmd-cli.cpp.o
[ 98%] Linking CXX executable ../../bin/llama-mtmd-cli
[ 98%] Built target llama-mtmd-cli
[ 99%] Building CXX object tools/cvector-generator/CMakeFiles/llama-cvector-generator.dir/cvector-generator.cpp.o
[ 99%] Linking CXX executable ../../bin/llama-cvector-generator
[ 99%] Built target llama-cvector-generator
[ 99%] Building CXX object tools/export-lora/CMakeFiles/llama-export-lora.dir/export-lora.cpp.o
[100%] Linking CXX executable ../../bin/llama-export-lora
[100%] Built target llama-export-lora
sanchari@ubuntults:~/llama.cpp$ ./build/bin/llama-bench -m ./gpt2-medium.gguf -p 0 -n 256 -t 1
| model           | size   | params | backend | threads | test          | t/s |
| -----           | -----: | -----: | -----: | -----: | -----:-----: | -----: |
| gpt2 0.4B F16  | 679.38 MiB | 354.82 M | BLAS    | 1        | tg256          | 14.66 ± 0.08 |
build: 9b17d74ab (7062)
sanchari@ubuntults:~/llama.cpp$ █

```

Explanation

Field	Value	Explanation
Model	gpt2 0.4B F16	Same GPT-2 Medium FP16 model used across all tasks.
Model Size	679.38 MiB	FP16 GGUF model size (unchanged).
Parameters	354.82 M	Same as Task 3.
Backend	BLAS	OpenBLAS acceleration enabled → faster matrix multiplications.
Threads	1	Still single-threaded, but BLAS optimizes GEMM operations.
Test	tg256	Same tg256 benchmark for fair comparison.
Throughput (t/s)	14.66 ± 0.08 tokens/sec	~10x faster than Task 3 due to BLAS-optimized matrix ops.

Task 6 – Floating-Point Performance Counters

In this task, you will collect and report all available floating-point related performance counters using the perf tool. This step will help understand the micro-architectural features of your CPU when running the LLM benchmark.

Steps.

- List down all the relevant floating point and memory traffic performance counters available on your machine for computing Operation

Intensity.

- Examples of counters to report (exact names may vary by CPU):
 - Floating Point Counters:
 - * fp arith inst retired.scalar single
 - * fp arith inst retired.scalar double
 - * fp arith inst retired.128b packed single
 - * fp arith inst retired.256b packed single
 - * fp arith inst retired.128b packed double
 - * fp arith inst retired.256b packed double
 - * fp arith inst retired.vector
 - Memory Counter:
 - * uncore imc free running/data total

Deliverables.

- Tabulate the counters with a short description of what each counter measures.

Solution:

Collected using:

```
perf list | grep -i "arith"
```

Outputs:

Floating-Point Counters

:

```
sanchari@ubuntults:~/llama.cpp$ perf list | grep -i "arith"
fp_arith_inst_retired.128b_packed_double
fp_arith_inst_retired.128b_packed_single
fp_arith_inst_retired.256b_packed_double
fp_arith_inst_retired.256b_packed_single
fp_arith_inst_retired.4_flops
fp_arith_inst_retired.512b_packed_double
fp_arith_inst_retired.512b_packed_single
fp_arith_inst_retired.8_flops
fp_arith_inst_retired.scalar
fp_arith_inst_retired.scalar_double
fp_arith_inst_retired.scalar_single
fp_arith_inst_retired.vector
    [Number of any Vector retired FP arithmetic instructions]
arith.divider_active
Error: failed to open tracing events directory
    [This metric approximates arithmetic floating-point (FP) scalar uops
    [This metric approximates arithmetic floating-point (FP) vector uops
    [This metric approximates arithmetic FP vector uops fraction the CPU
    [This metric approximates arithmetic FP vector uops fraction the CPU
    [This metric approximates arithmetic FP vector uops fraction the CPU
tma_info_core_fp_arith_utilization
    [This metric approximates arithmetic floating-point (FP) scalar uops
    [This metric approximates arithmetic floating-point (FP) vector uops
    [This metric approximates arithmetic FP vector uops fraction the CPU
    [This metric approximates arithmetic FP vector uops fraction the CPU
    [This metric approximates arithmetic FP vector uops fraction the CPU
tma_info_core_fp_arith_utilization
tma_info_inst_mix_iparith
    [Instructions per FP Arithmetic instruction (lower number means higher
tma_info_inst_mix_iparith_avx128
    [Instructions per FP Arithmetic AVX/SSE 128-bit instruction (lower
tma_info_inst_mix_iparith_avx256
    [Instructions per FP Arithmetic AVX* 256-bit instruction (lower number
tma_info_inst_mix_iparith_avx512
    [Instructions per FP Arithmetic AVX 512-bit instruction (lower number
tma_info_inst_mix_iparith_scalar_dp
    [Instructions per FP Arithmetic Scalar Double-Precision instruction
tma_info_inst_mix_iparith_scalar_sp
    [Instructions per FP Arithmetic Scalar Single-Precision instruction
tma_info_inst_mix_iparith_scalar_dp
    [Instructions per FP Arithmetic Scalar Double-Precision instruction
tma_info_inst_mix_iparith_scalar_sp
    [Instructions per FP Arithmetic Scalar Single-Precision instruction
tma_info_inst_mix_iparith_avx128
    [Instructions per FP Arithmetic AVX/SSE 128-bit instruction (lower
tma_info_inst_mix_iparith_avx256
    [Instructions per FP Arithmetic AVX* 256-bit instruction (lower number
```

```

tma_info_inst_mix_iparith_avx512
    [Instructions per FP Arithmetic AVX 512-bit instruction (lower number
tma_fp_arith
    [This metric represents overall arithmetic floating-point (FP)
tma_info_core_fp_arith_utilization
tma_info_inst_mix_iparith
    [Instructions per FP Arithmetic instruction (lower number means higher
tma_info_inst_mix_iparith_avx128
    [Instructions per FP Arithmetic AVX/SSE 128-bit instruction (lower
tma_info_inst_mix_iparith_avx256
    [Instructions per FP Arithmetic AVX* 256-bit instruction (lower number
tma_info_inst_mix_iparith_avx512
    [Instructions per FP Arithmetic AVX 512-bit instruction (lower number
tma_info_inst_mix_iparith_scalar_dp
    [Instructions per FP Arithmetic Scalar Double-Precision instruction
tma_info_inst_mix_iparith_scalar_sp
    [Instructions per FP Arithmetic Scalar Single-Precision instruction
tma_fp_arith
    [This metric represents overall arithmetic floating-point (FP)
    [This metric approximates arithmetic floating-point (FP) scalar uops
    [This metric approximates arithmetic floating-point (FP) vector uops
    [This metric approximates arithmetic FP vector uops fraction the CPU
    [This metric approximates arithmetic FP vector uops fraction the CPU
    [This metric approximates arithmetic FP vector uops fraction the CPU
tma_fp_arith
    [This metric represents overall arithmetic floating-point (FP)
    [This metric approximates arithmetic floating-point (FP) scalar uops
    [This metric approximates arithmetic floating-point (FP) vector uops
    [This metric approximates arithmetic FP vector uops fraction the CPU
    [This metric approximates arithmetic FP vector uops fraction the CPU
    [This metric approximates arithmetic FP vector uops fraction the CPU
tma_fp_arith_group: [Metrics contributing to tma_fp_arith category]
    [This metric approximates arithmetic floating-point (FP) scalar uops
    [This metric approximates arithmetic floating-point (FP) vector uops
    [This metric approximates arithmetic FP vector uops fraction the CPU
    [This metric approximates arithmetic FP vector uops fraction the CPU
    [This metric approximates arithmetic FP vector uops fraction the CPU
    [This metric approximates arithmetic floating-point (FP) scalar uops
    [This metric approximates arithmetic floating-point (FP) vector uops
    [This metric approximates arithmetic FP vector uops fraction the CPU
    [This metric approximates arithmetic FP vector uops fraction the CPU
    [This metric approximates arithmetic FP vector uops fraction the CPU
tma_fp_arith
    [This metric represents overall arithmetic floating-point (FP)
sanchari@ubuntults:~/llama.cpp$ █

```

Memory Counters :

perf list | grep -i "memory"

```
[sanchari@ubuntults: ~/llama.cpp]
sanchari@ubuntults:~/llama.cpp$ perf list | grep -i "memory"
    directly written to memory and not allocated in L3. Clean lines may
    [All retired memory instructions Supports address when precise (Precise
    [Any memory transaction that reached the SQ]

memory:
    [Number of times an HLE execution aborted due to various memory events
    [Number of times an HLE execution aborted due to incompatible memory
machine_clears.memory_ordering
    [Counts the number of machine clears due to memory order conflicts Spec
    [Number of times an RTM execution aborted due to various memory events
    [Number of times an RTM execution aborted due to incompatible memory
memory_disambiguation.history_reset
    [MEMORY_DISAMBIGUATION.HISTORY_RESET]
    [Cycles while memory subsystem has an outstanding load]
    [Execution stalls while memory subsystem has an outstanding load]
        clear event for this thread (e.g. misprediction or memory nuke)]
        misprediction or memory nuke)
virtual memory:
Error: failed to open tracing events directory
    [This metric represents fraction of slots where Core non-memory issues
tma_info_memory_12mpki
tma_memory_bound
    [This metric represents fraction of slots the Memory subsystem within
tma_info_memory_fb_hpki
tma_info_memory_11mpki
tma_info_memory_11mpki_load
tma_info_memory_12hpki_all
tma_info_memory_12hpki_load
tma_info_memory_12mpki
tma_info_memory_12mpki_all
tma_info_memory_12mpki_load
tma_info_memory_13mpki
    [This metric represents fraction of slots where Core non-memory issues
    [This metric estimates fraction of cycles while the memory subsystem
tma_info_memory_tlb_code_stlb_mpki
    [Average external Memory Bandwidth Use for reads and writes [GB / sec]]
tma_info_memory_core_12_evictions_nonsilent_pkis
tma_info_memory_core_12_evictions_silent_pkis
    the evicted lines are dropped (no writeback to L3 or memory)
tma_info_bottleneck_memory_bandwidth
    [Total pipeline cost of (external) Memory Bandwidth related bottlenecks]
tma_info_bottleneck_memory_data_tlbs
    [Total pipeline cost of Memory Address Translation related bottlenecks]
tma_info_bottleneck_memory_latency
    [Total pipeline cost of Memory Latency related bottlenecks (external
    memory and off-core caches)]
```

```

    memory and off-core caches)]
tma_info_memory_core_l1d_cache_fill_bw
tma_info_memory_core_12_cache_fill_bw
tma_info_memory_core_12_evictions_nonsilent_pk
tma_info_memory_core_12_evictions_silent_pk
    the evicted lines are dropped (no writeback to L3 or memory)
tma_info_memory_core_13_cache_access_bw
tma_info_memory_core_13_cache_fill_bw
tma_info_memory_fb_hpki
tma_info_memory_l1mpki
tma_info_memory_l1mpki_load
tma_info_memory_12hpki_all
tma_info_memory_12hpki_load
tma_info_memory_12mpki
tma_info_memory_12mpki_all
tma_info_memory_12mpki_load
tma_info_memory_13mpki
tma_info_memory_load_miss_real_latency
tma_info_memory_mlp
    [Memory-Level-Parallelism (average number of L1 miss demand load when
tma_info_memory_thread_l1d_cache_fill_bw_lt
tma_info_memory_thread_12_cache_fill_bw_lt
tma_info_memory_thread_13_cache_access_bw_lt
tma_info_memory_thread_13_cache_fill_bw_lt
tma_info_memory_tlb_load_stlb_mpki
tma_info_memory_tlb_page_walks_utilization
tma_info_memory_tlb_store_stlb_mpki
    [Average external Memory Bandwidth Use for reads and writes [GB / sec]]
    [Average latency of data read request to external DRAM memory [in
    [Average number of parallel data read requests to external memory]
    [Average latency of data read request to external memory (in
MemoryBW: [Grouping from Top-down Microarchitecture Analysis Metrics spreadsheet]
    unavailability limited additional L1D miss memory access requests to
tma_info_bottleneck_memory_bandwidth
    [Total pipeline cost of (external) Memory Bandwidth related bottlenecks]
tma_info_memory_core_l1d_cache_fill_bw
tma_info_memory_core_12_cache_fill_bw
tma_info_memory_core_13_cache_access_bw
tma_info_memory_core_13_cache_fill_bw
tma_info_memory_mlp
    [Memory-Level-Parallelism (average number of L1 miss demand load when
tma_info_memory_thread_l1d_cache_fill_bw_lt
tma_info_memory_thread_12_cache_fill_bw_lt
tma_info_memory_thread_13_cache_access_bw_lt
tma_info_memory_thread_13_cache_fill_bw_lt
    [Average external Memory Bandwidth Use for reads and writes [GB / sec]]
    [Average number of parallel data read requests to external memory]
    was likely hurt due to approaching bandwidth limits of external memory

```

```

MemoryBound: [Grouping from Top-down Microarchitecture Analysis Metrics spreadsheet]
    external memory (DRAM) by loads
    tma_info_memory_load_miss_real_latency
    tma_info_memory_mlp
        [Memory-Level-Parallelism (average number of L1 miss demand load when
            memory accesses; RFO store issue a read-for-ownership request before
MemoryLat: [Grouping from Top-down Microarchitecture Analysis Metrics spreadsheet]
    tma_info_bottleneck_memory_latency
        [Total pipeline cost of Memory Latency related bottlenecks (external
            memory and off-core caches)]
    tma_info_memory_load_miss_real_latency
        [Average latency of data read request to external DRAM memory [in
            [Average latency of data read request to external memory (in
                likely hurt due to latency from external memory (DRAM)]]
MemoryTLB: [Grouping from Top-down Microarchitecture Analysis Metrics spreadsheet]
    tma_info_bottleneck_memory_data_tlbs
        [Total pipeline cost of Memory Address Translation related bottlenecks
    tma_info_memory_tlb_code_stlb_mpki
    tma_info_memory_tlb_load_stlb_mpki
    tma_info_memory_tlb_page_walks_utilization
    tma_info_memory_tlb_store_stlb_mpki
Memory_BW: [Grouping from Top-down Microarchitecture Analysis Metrics spreadsheet]
    tma_info_memory_oro_data_12_mlp
    tma_info_memory_oro_load_12_mlp
Memory_Lat: [Grouping from Top-down Microarchitecture Analysis Metrics spreadsheet]
    tma_info_memory_oro_load_12_miss_latency
        that are requesting memory from the CPU]
        that are writing memory to the CPU]
        read miss (read memory access) in nano seconds]
        read miss (read memory access) addressed to local memory in nano
        read miss (read memory access) addressed to remote memory in nano
    llc_miss_local_memory_bandwidth_read
        (LLC) and go to local memory]
    llc_miss_local_memory_bandwidth_write
        (LLC) and go to local memory]
    llc_miss_remote_memory_bandwidth_read
        (LLC) and go to remote memory]
        [The ratio of number of completed memory load instructions to the total
    memory_bandwidth_read
        [DDR memory read bandwidth (MB/sec)]
    memory_bandwidth_total
        [DDR memory bandwidth (MB/sec)]
    memory_bandwidth_write
        [DDR memory write bandwidth (MB/sec)]

```

```

memory_bandwidth_write
    [DDR memory write bandwidth (MB/sec)]
[Memory read that miss the last level cache (LLC) addressed to local
DRAM as a percentage of total memory read accesses, does not include
[Memory reads that miss the last level cache (LLC) addressed to remote
DRAM as a percentage of total memory read accesses, does not include
[Cycles Memory is in self refresh power mode]
[The ratio of number of completed memory store instructions to the
[This metric estimates fraction of cycles while the memory subsystem
[This metric estimates fraction of cycles while the memory subsystem
tma_info_bottleneck_memory_bandwidth
    [Total pipeline cost of (external) Memory Bandwidth related bottlenecks]
tma_info_bottleneck_memory_data_tlbs
    [Total pipeline cost of Memory Address Translation related bottlenecks
tma_info_bottleneck_memory_latency
    [Total pipeline cost of Memory Latency related bottlenecks (external
        memory and off-core caches)]
tma_info_memory_core_13_cache_access_bw
tma_info_memory_12mpki_all
tma_info_memory_oro_data_12_mlp
tma_info_memory_oro_load_12_miss_latency
tma_info_memory_oro_load_12_mlp
tma_info_memory_thread_13_cache_access_bw_lt
    was likely hurt due to approaching bandwidth limits of external memory
    likely hurt due to latency from external memory (DRAM)
    [This metric estimates fraction of cycles while the memory subsystem
tma_memory_operations
    memory operations -- uops for memory load or store accesses]
tma_info_memory_core_12_evictions_nonsilent_pk
tma_info_memory_core_12_evictions_silent_pk
    the evicted lines are dropped (no writeback to L3 or memory)
    [Average latency of data read request to external DRAM memory [in
    [This metric estimates fraction of cycles while the memory subsystem
        was handling loads from local memory]
    [This metric estimates fraction of cycles while the memory subsystem
    [This metric estimates fraction of cycles while the memory subsystem
        was handling loads from remote memory]
    [This metric estimates fraction of cycles while the memory subsystem
    [This metric estimates fraction of cycles while the memory subsystem
    [This metric estimates fraction of cycles while the memory subsystem
        was handling loads from remote memory]
    [Average external Memory Bandwidth Use for reads and writes [GB / sec]]
    [Average latency of data read request to external DRAM memory [in
    [Average number of parallel data read requests to external memory]
    [Average latency of data read request to external memory (in
    [This metric represents fraction of slots where Core non-memory issues
tma_memory_bound
    [This metric represents fraction of slots the Memory subsystem within

```

```
    external memory (DRAM) by loads]
    memory accesses; RFO store issue a read-for-ownership request before
    [This metric represents fraction of slots where Core non-memory issues
tma_memory_bound
    [This metric represents fraction of slots the Memory subsystem within
     external memory (DRAM) by loads]
tma_memory_operations
    memory operations -- uops for memory load or store accesses]
    memory accesses; RFO store issue a read-for-ownership request before
    [This metric estimates how often memory load accesses were aliased by
    [This metric estimates fraction of cycles while the memory subsystem
    [This metric estimates fraction of cycles while the memory subsystem
unavailability limited additional L1D miss memory access requests to
was likely hurt due to approaching bandwidth limits of external memory
likely hurt due to latency from external memory (DRAM)]
[This metric estimates fraction of cycles handling memory load split
[This metric roughly estimates fraction of cycles when the memory
[This metric estimates fraction of cycles while the memory subsystem
was handling loads from local memory]
[This metric estimates fraction of cycles while the memory subsystem
[This metric estimates fraction of cycles while the memory subsystem
was handling loads from remote memory]
[This metric represents fraction of slots where Core non-memory issues
tma_memory_bound
    [This metric represents fraction of slots the Memory subsystem within
     external memory (DRAM) by loads]
tma_memory_operations
    memory operations -- uops for memory load or store accesses]
    memory accesses; RFO store issue a read-for-ownership request before
    [This metric estimates how often memory load accesses were aliased by
    [This metric estimates fraction of cycles while the memory subsystem
    [This metric estimates fraction of cycles while the memory subsystem
unavailability limited additional L1D miss memory access requests to
was likely hurt due to approaching bandwidth limits of external memory
likely hurt due to latency from external memory (DRAM)]
[This metric estimates fraction of cycles handling memory load split
[This metric roughly estimates fraction of cycles when the memory
[This metric estimates fraction of cycles while the memory subsystem
was handling loads from local memory]
[This metric estimates fraction of cycles while the memory subsystem
[This metric estimates fraction of cycles while the memory subsystem
was handling loads from remote memory]
[This metric represents fraction of slots where Core non-memory issues
tma_memory_bound
    [This metric represents fraction of slots the Memory subsystem within
     was likely hurt due to approaching bandwidth limits of external memory
likely hurt due to latency from external memory (DRAM)]
unavailability limited additional L1D miss memory access requests to
tma_info_bottleneck_memory_bandwidth
```

```

[This metric estimates fraction of cycles while the memory subsystem
unavailability limited additional L1D miss memory access requests to
was likely hurt due to approaching bandwidth limits of external memory
likely hurt due to latency from external memory (DRAM)]
[This metric estimates fraction of cycles handling memory load split
[This metric roughly estimates fraction of cycles when the memory
[This metric estimates fraction of cycles while the memory subsystem
was handling loads from local memory]
[This metric estimates fraction of cycles while the memory subsystem
[This metric estimates fraction of cycles while the memory subsystem
was handling loads from remote memory]
[This metric represents fraction of slots where Core non-memory issues
tma_memory_bound
    [This metric represents fraction of slots the Memory subsystem within
    was likely hurt due to approaching bandwidth limits of external memory
    likely hurt due to latency from external memory (DRAM)]
    unavailability limited additional L1D miss memory access requests to
tma_info_bottleneck_memory_bandwidth
    [Total pipeline cost of (external) Memory Bandwidth related bottlenecks]
    [Average external Memory Bandwidth Use for reads and writes [GB / sec]]
    was likely hurt due to approaching bandwidth limits of external memory
tma_info_bottleneck_memory_latency
    [Total pipeline cost of Memory Latency related bottlenecks (external
    memory and off-core caches)]
    likely hurt due to latency from external memory (DRAM)
    unavailability limited additional L1D miss memory access requests to
    unavailability limited additional L1D miss memory access requests to
    [This metric estimates fraction of cycles while the memory subsystem
    [This metric estimates fraction of cycles while the memory subsystem
    [This metric estimates fraction of cycles while the memory subsystem
tma_info_bottleneck_memory_data_tlbs
    [Total pipeline cost of Memory Address Translation related bottlenecks
    [This metric estimates how often memory load accesses were aliased by
    unavailability limited additional L1D miss memory access requests to
    [This metric estimates fraction of cycles handling memory load split
    [This metric roughly estimates fraction of cycles when the memory
    [This metric estimates fraction of cycles while the memory subsystem
    [This metric estimates fraction of cycles while the memory subsystem
tma_memory_operations
    memory operations -- uops for memory load or store accesses]
    [This metric estimates fraction of cycles while the memory subsystem
    was handling loads from local memory]
    [This metric estimates fraction of cycles while the memory subsystem
    [This metric estimates fraction of cycles while the memory subsystem
    was handling loads from remote memory]
tma_memory_bound_group: [Metrics contributing to tma_memory_bound category]
    external memory (DRAM) by loads
        memory accesses; RFO store issue a read-for-ownership request before
sanchari@ubuntults:~/llama.cpp$ 

```

These counters are used for computing **Operation Intensity (OI)** and locating performance bounds in the Roofline model.

Counter / Metric Name	Description
fp_arith_inst_retired.scalar_single	Retired scalar single-precision (32-bit) FP arithmetic instructions.
fp_arith_inst_retired.scalar_double	Retired scalar double-precision (64-bit) FP arithmetic instructions.
fp_arith_inst_retired.scalar	Total number of scalar FP arithmetic instructions retired (all precisions).
fp_arith_inst_retired.vector	Total number of vector (SIMD) FP arithmetic instructions retired.
fp_arith_inst_retired.128b_packed_single	Retired 128-bit packed single-precision FP instructions (SSE).

fp_arith_inst_retired.128b_packed_double	Retired 128-bit packed double-precision FP instructions (SSE2).
fp_arith_inst_retired.256b_packed_single	Retired 256-bit packed single-precision FP instructions (AVX).
fp_arith_inst_retired.256b_packed_double	Retired 256-bit packed double-precision FP instructions (AVX).
fp_arith_inst_retired.512b_packed_single	Retired 512-bit packed single-precision FP instructions (AVX-512).
fp_arith_inst_retired.512b_packed_double	Retired 512-bit packed double-precision FP instructions (AVX-512).
fp_arith_inst_retired.4_flops	Represents 4 floating-point operations per instruction (e.g., SSE).
fp_arith_inst_retired.8_flops	Represents 8 floating-point operations per instruction (e.g., AVX-512).
fp_assist.any	Cycles with any FP assist or exception handling (e.g., denormals, div-by-zero).
tma_fp_arith	Overall metric representing total arithmetic FP operations executed.
tma_fp_scalar	Approximates arithmetic scalar FP uops utilization.
tma_fp_vector	Approximates arithmetic vector FP uops utilization (all widths).
tma_fp_vector_128b	Fraction of CPU pipeline executing 128-bit FP vector operations.
tma_fp_vector_256b	Fraction of CPU pipeline executing 256-bit FP vector operations.
tma_fp_vector_512b	Fraction of CPU pipeline executing 512-bit FP vector operations.
tma_info_core_fp_arith_utilization	FP arithmetic utilization efficiency — lower value means higher FP usage efficiency.

Counter Name	Description
memory_bandwidth_read	DDR memory read bandwidth (MB/sec).
memory_bandwidth_write	DDR memory write bandwidth (MB/sec).
memory_bandwidth_total	Total DDR memory bandwidth (MB/sec) = read + write.
llc_miss_local_memory_bandwidth_read	Bandwidth of memory reads that miss LLC and go to local DRAM .
llc_miss_remote_memory_bandwidth_read	Bandwidth of memory reads that miss LLC and go to remote DRAM (NUMA).
tma_memory_bound	Fraction of CPU pipeline slots stalled due to memory subsystem bottlenecks.
tma_info_memory_mlp	Average number of parallel L1-miss demand loads (Memory-Level Parallelism).
tma_info_memory_load_miss_real_latency	Average latency (in cycles) for data read requests to external DRAM.
tma_info_memory_core_l3_cache_fill_bw	Bandwidth used for L3 cache line fills from memory.
tma_info_memory_core_l2_cache_fill_bw	Bandwidth used for L2 cache line fills .

tma_info_memory_core_l1d_cache_fill_bw	Bandwidth used for L1 data cache fills .
tma_info_bottleneck_memory_bandwidth	Pipeline cost due to memory bandwidth saturation .
tma_info_bottleneck_memory_latency	Pipeline cost due to memory latency from DRAM or off-core caches.
uncore_imc_free_running/data_total	Integrated Memory Controller (IMC) counter – measures total memory traffic (reads + writes).

Event / Metric Name	Description / Meaning
fp_arith_inst_retired.128b_packed_double	Retired 128-bit packed double-precision FP instructions
fp_arith_inst_retired.128b_packed_single	Retired 128-bit packed single-precision FP instructions
fp_arith_inst_retired.256b_packed_double	Retired 256-bit packed double-precision FP instructions
fp_arith_inst_retired.256b_packed_single	Retired 256-bit packed single-precision FP instructions
fp_arith_inst_retired.512b_packed_double	Retired 512-bit packed double-precision FP instructions
fp_arith_inst_retired.512b_packed_single	Retired 512-bit packed single-precision FP instructions
fp_arith_inst_retired.4_flops	Four floating-point operations (FLOPs) per instruction retired
fp_arith_inst_retired.8_flops	Eight floating-point operations (FLOPs) per instruction retired
fp_arith_inst_retired.scalar	Scalar FP arithmetic instructions retired
fp_arith_inst_retired.scalar_double	Scalar double-precision FP arithmetic instructions retired
fp_arith_inst_retired.scalar_single	Scalar single-precision FP arithmetic instructions retired
fp_arith_inst_retired.vector	Number of vector FP arithmetic instructions retired

Event / Metric Name	Description / Meaning
arith.divider_active	Fraction of cycles where the FP divider unit is active (arithmetic operations in progress)
tma_fp_arith	Overall floating-point arithmetic operations utilization metric
tma_info_core_fp_arith_utilization	Measures utilization of FP arithmetic instructions (lower = higher utilization)

Metric Name	Meaning
tma_info_inst_mix_iparith	Instructions per FP arithmetic instruction (lower = higher FP activity)
tma_info_inst_mix_iparith_avx128	Instructions per 128-bit AVX/SSE FP arithmetic instruction
tma_info_inst_mix_iparith_avx256	Instructions per 256-bit AVX FP arithmetic instruction

tma_info_inst_mix_iparith_avx512	Instructions per 512-bit AVX FP arithmetic instruction
tma_info_inst_mix_iparith_scalar_dp	Instructions per scalar double-precision FP arithmetic instruction
tma_info_inst_mix_iparith_scalar_sp	Instructions per scalar single-precision FP arithmetic instruction

Group / Metric	Description / Role
tma_fp_arith_group	Aggregate metric group for FP arithmetic instructions
tma_fp_arith	Represents total FP arithmetic activity (scalar + vector)
tma_fp_scalar	Scalar floating-point uops retired
tma_fp_vector	Vector floating-point uops retired
tma_fp_vector_128b	128-bit FP vector uops (SSE/AVX)
tma_fp_vector_256b	256-bit FP vector uops (AVX2)
tma_fp_vector_512b	512-bit FP vector uops (AVX-512)

Category	Representative Events
Scalar FP	fp_arith_inst_retired.scalar*, tma_fp_scalar, tma_info_inst_mix_iparith_scalar_*
Vector FP	fp_arith_inst_retired.{128b,256b,512b}_*, tma_fp_vector_*
Overall Arithmetic Load	tma_fp_arith, tma_info_core_fp_arith_utilization, arith.divider_active

Task 7 – Performance Counters and Roofline Analysis

In this task, you will collect performance counters for all build variants of llama.cpp (Tasks 3, 4, and 5) and use them to derive Operation Intensity (OI) and perform a Roofline analysis.

Deliverables.

- Submit the performance counter values collected using perf for each variant.
- Provide the Operation Intensity derivation for each case.
- Submit peak memory bandwidth and peak compute capacity values of your system.
- Submit the Roofline analysis plot overlaying all three variants.
- Provide a short commentary on where each variant sits (memory-bound vs compute-bound).

Solution:

Task3 Details:

Field	Value	Explanation
Model	gpt2 0.4B F16	GPT-2 Medium in FP16 GGUF format (\approx 354M parameters).
Model Size	679.38 MiB	FP16 GGUF file loaded by llama.cpp.
Parameters	354.82 M	Total number of FP16 weights in the model.
Backend	CPU (Naive Build)	No SIMD, no MKL, no OpenBLAS
Threads	1	Task 3 requires <i>single-thread execution</i> to measure true naive performance.

Field	Value	Explanation
Test	tg256	Standard 256-token generation benchmark used by llama.cpp.
Throughput (t/s)	1.40 ± 0.01 tokens/sec	Extremely low speed due to no parallelism or CPU vectorization.

Task4 Details:

Field	Value	Explanation
Model	gpt2 0.4B F16	GPT-2 Medium model in FP16 GGUF format (~354M parameters).
Model Size	679.38 MiB	Size of the loaded FP16 model file.
Parameters	354.82 M	Number of neural network weights in GPT-2 Medium.
Backend	CPU	Naive llama.cpp build—no BLAS, no MKL, no SIMD acceleration.
Threads	1	Task 4 is evaluated single-threaded to measure baseline performance.
Test	tg256	Benchmarks token-generation performance for 256 tokens.
Throughput (t/s)	1.40 ± 0.00 tokens/sec	Slow performance; reflects absence of SIMD or BLAS optimizations.

Task5 Details:

Field	Value	Explanation
Model	gpt2 0.4B F16	Same GPT-2 Medium FP16 model used across all tasks.
Model Size	679.38 MiB	FP16 GGUF model size (unchanged).
Parameters	354.82 M	Same as Task 3.
Backend	BLAS	OpenBLAS acceleration enabled → faster matrix multiplications.
Threads	1	Still single-threaded, but BLAS optimizes GEMM operations.
Test	tg256	Same tg256 benchmark for fair comparison.
Throughput (t/s)	14.66 ± 0.08 tokens/sec	~10x faster than Task 3 due to BLAS-optimized matrix ops.

Snapshots:

Collecting Performance Counters for Task3, Task4 and Task5:

How to run:

```
perf stat -o perf_variant3_data.txt -e fp_arith_inst_retired.scalar -e fp_arith_inst_retired.4_flops -e fp_arith_inst_retired.8_flops -e fp_arith_inst_retired.256b_packed_single -e cycles -e cache-misses -e cache-references -- ./build-variant-generic/bin/llama-bench -m gpt2-medium.gguf -p 0 -n 256 -t 1
```

-o perf_variant3_data.txt :Writes perf output to a file.

FP Instruction Counters : These counters measure floating-point arithmetic instructions:

fp_arith_inst_retired.scalar :FP instructions that do 1 FLOP

fp_arith_inst_retired.4_flops : Vector instructions performing 4 FLOPs

fp_arith_inst_retired.8_flop : Vector instructions performing 8 FLOPs : AVX2 (256-bit registers with 8 lanes)

fp_arith_inst_retired.256b_packed_single :FP instructions operating on full 256-bit registers

Cycles:Tells how many CPU cycles the program consumed

cache-misses : Counts how many cache lines were missed.

cache-references :Total cache accesses.

Running the Benchmark ./llama-bench -m gpt2-medium.gguf -p 0 -n 256 -t 1

Runs the actual GPT-2 inference:

- -m → model file
- -p 0 → prompt index 0
- -n 256 → generate 256 tokens
- -t 1 → 1 inference thread

Task3:

```
sanchari@ubuntults:~/llama.cpp$ ./llama.cpp$ cmake -B build -DGML_CPU_GENERIC=ON -DGML_NATIVE=OFF \
-dGML_AVX=OFF -DGML_AVX2=OFF -DGML_AVX512=OFF \
-dGML_SSE42=OFF -DGML_F16C=OFF -DGML_FMA=OFF
cmake --build build --config Release
CMAKE BUILD TYPE=Release
-- ccache found, compilation results will be cached. Disable with GGML_CCACHE=OFF.
-- CMAKE SYSTEM PROCESSOR: x86_64
-- GGML_SYSTEM_ARCH: x86
-- Including CPU backend
-- x86 detected
-- Adding CPU backend variant ggml-cpu:
-- BLAS found, Libraries: /opt/intel/oneapi/mkl/2025.3/lib/libmkl_intel_lp64.so;/opt/intel/oneapi/mkl/2025.3/lib/libmkl_intel_thread.so;/opt/intel/oneapi/mkl/intel/oneapi/compiler/2025.3/lib/libimc5.so;-lm;-ldl
-- BLAS found, Includes: /opt/intel/oneapi/mkl/2025.3/lib/pkgconfig/.../include
-- Including BLAS backend
-- ggml version: 0.9.4
-- ggml commit: 9b17d74ab
-- Configuring done (0.6s)
-- Generating done (0.5s)
-- Build files have been written to: /home/sanchari/llama.cpp/build
[ 3%] Built target ggml-base
[ 3%] Built target ggml-bias
[ 4%] Building C object ggml/src/CMakeFiles/ggml-cpu.dir/ggml-cpu/ggml-cpu.c.o
[ 4%] Building CXX object ggml/src/CMakeFiles/ggml-cpu.dir/ggml-cpu/ggml-cpu.cpp.o
[ 4%] Building CXX object ggml/src/CMakeFiles/ggml-cpu.dir/ggml-cpu/repack.cpp.o
[ 4%] Building CXX object ggml/src/CMakeFiles/ggml-cpu.dir/ggml-cpu/hbm.cpp.o
[ 5%] Building C object ggml/src/CMakeFiles/ggml-cpu.dir/ggml-cpu/quants.c.o
[ 5%] Building CXX object ggml/src/CMakeFiles/ggml-cpu.dir/ggml-cpu/traits.cpp.o
[ 5%] Building CXX object ggml/src/CMakeFiles/ggml-cpu.dir/ggml-cpu/amx/amx.cpp.o
[ 6%] Building CXX object ggml/src/CMakeFiles/ggml-cpu.dir/ggml-cpu/amx/mmq.cpp.o
[ 6%] Building CXX object ggml/src/CMakeFiles/ggml-cpu.dir/ggml-cpu/binary-ops.cpp.o
[ 6%] Building CXX object ggml/src/CMakeFiles/ggml-cpu.dir/ggml-cpu/unary-ops.cpp.o
[ 6%] Building CXX object ggml/src/CMakeFiles/ggml-cpu.dir/ggml-cpu/vec.cpp.o
[ 7%] Building CXX object ggml/src/CMakeFiles/ggml-cpu.dir/ggml-cpu/ops.cpp.o
[ 7%] Building CXX object ggml/src/CMakeFiles/ggml-cpu.dir/ggml-cpu/arch/x86/quants.c.o
[ 8%] Building CXX object ggml/src/CMakeFiles/ggml-cpu.dir/ggml-cpu/arch/x86/repack.cpp.o
[ 8%] Linking CXX shared library ../../bin/libggml-cpu.so
[ 8%] Built target ggml-cpu
[ 8%] Linking CXX shared library ../../bin/libggml.so
[ 8%] Built target ggml
[ 8%] Linking CXX shared library ../../bin/libllama.so
[ 43%] Built target llama
[ 43%] Built target build_info
[ 47%] Built target common
[ 48%] Built target cpp-httplib
[ 48%] Linking CXX executable ../../bin/test-tokenizer-0
[ 50%] Built target llama-qdqd
[ 90%] Linking CXX executable ../../bin/llama-batched-bench
[ 90%] Built target llama-batched-bench
[ 90%] Linking CXX executable ../../bin/llama-gguf-split
[ 91%] Built target llama-gguf-split
[ 92%] Linking CXX executable ../../bin/llama-imatrix
[ 92%] Built target llama-imatrix
[ 93%] Linking CXX executable ../../bin/llama-bench
[ 93%] Built target llama-bench
[ 93%] Linking CXX executable ../../bin/llama-cli
[ 93%] Built target llama-cli
[ 93%] Linking CXX executable ../../bin/llama-perplexity
[ 93%] Built target llama-perplexity
[ 94%] Linking CXX executable ../../bin/llama-quantize
[ 94%] Built target llama-quantize
[ 94%] Linking CXX executable ../../bin/llama-server
[ 95%] Built target llama-server
[ 95%] Linking CXX executable ../../bin/llama-run
[ 95%] Built target llama-run
[ 96%] Linking CXX executable ../../bin/llama-tokenize
[ 96%] Built target llama-tokenize
[ 97%] Linking CXX executable ../../bin/llama-tts
[ 97%] Built target llama-tts
[ 97%] Built target llama-llava-cli
[ 98%] Built target llama-gemm3-cli
[ 98%] Built target llama-minicpmv-cli
[ 98%] Built target llama-qwen2v1-cli
[ 99%] Linking CXX executable ../../bin/llama-mtmd-cli
[ 99%] Built target llama-mtmd-cli
[ 99%] Linking CXX executable ../../bin/llama-cvector-generator
[ 99%] Built target llama-cvector-generator
[100%] Linking CXX executable ../../bin/llama-export-lora
sanchari@ubuntults:~/llama.cpp$ sanchari@ubuntults:~/llama.cpp$ sanchari@ubuntults:~/llama.cpp$ sanchari@ubuntults:~/llama.cpp$ sanchari@ubuntults:~/llama.cpp$ sanchari@ubuntults:~/llama.cpp$ sanchari@ubuntults:~/llama.cpp$ sanchari@ubuntults:~/llama.cpp$ sanchari@ubuntults:~/llama.cpp$ perf stat -o perf.variant3.data.txt -e fp_arith_inst_retired.scalar -e fp_arith_inst_retired.4_flops -e fp_arith_inst_retired.8_flops -e fp_arith_inst_retire
d.256_packed_single -e cycles -e cache-misses -e cache-references -- ./build-variant-generic/bin/llama-bench -m gpt2-medium.gguf -p 0 -n 256 -t 1
| model | size | params | backend | threads | test | t/s |
| ----- | ----- | ----- | ----- | ----- | ----- | ----- |
| gpt2_0.48_F16 | 679.38 MB | 354.82 M | CPU | 1 | tg256 | 1.39 ± 0.01 |
build: 9b17d74ab (7062)
sanchari@ubuntults:~/llama.cpp$
```

```
vi perf_variant3_data.txt
```

```
# started on Tue Nov 18 05:56:45 2025

Performance counter stats for './build-variant-generic/bin/llama-bench -m gpt2-medium.gguf -p 0 -n 256 -t 1':

  938,074,086,809    fp_arith_inst_retired.scalar
  970,362,624    fp_arith_inst_retired.4_flops
  970,362,624    fp_arith_inst_retired.8_flops
      0    fp_arith_inst_retired.256b_packed_single
1,913,320,659,935    cycles
14,118,957,840    cache-misses          #  88.42% of all cache refs
15,967,598,832    cache-references

  921.169631290 seconds time elapsed

  920.954707000 seconds user
  0.155982000 seconds sys
```

Task 4:

```
sanchari@ubuntults:~/llama.cpp$ ./llama.cpp$ cmake -B build
CMAKE_BUILD_TYPE=release
-- ccache found, compilation results will be cached. Disable with GGML_CCACHE=OFF.
-- CMAKE_SYSTEM_PROCESSOR: x86_64
-- GGML_SYSTEM_ARCH: x86
-- Including CPU backend
-- x86 detected
-- Adding CPU backend variant ggml-cpu:
-- BLAS found, Libraries: /opt/intel/oneapi/mkl/2025.3/lib/libmkl_intel_lp64.so;/opt/intel/oneapi/mkl/2025.3/lib/libmkl_intel_thread.so;/opt/intel/oneapi/mkl/2025.3/lib/libmkl_core.so;-lm;-ldl
-- BLAS found, Includes: /opt/intel/oneapi/mkl/2025.3/lib/pkgconfig/../../include
-- Including BLAS backend
-- ggml version: 0.9.4
-- ggml commit: 9b17d74ab
-- Configuring done (0.6s)
-- Generating done (0.5s)
-- Build files have been written to: /home/sanchari/llama.cpp/build
sanchari@ubuntults:~/llama.cpp$ cmake --build build --config Release
[  3%] Built target ggml-base
[  8%] Built target ggml-bias
[  8%] Built target ggml-cpu
[  8%] Built target ggml
[ 43%] Built target llama
[ 43%] Built target build_info
[ 47%] Built target common
[ 48%] Built target cpp-httplib
[ 48%] Built target test-tokenizer-0
[ 49%] Built target test-sampling
[ 50%] Built target test-grammar-parser
[ 51%] Built target test-grammar-integration
[ 51%] Built target test-llama-grammar
[ 52%] Built target test-chat
[ 77%] Built target shal
[ 77%] Built target llama-gguf-hash
[ 78%] Built target llama-gguf
[ 78%] Built target llama-lookahead
[ 79%] Built target llama-lookup
[ 79%] Built target llama-lookup-create
[ 80%] Built target llama-lookup-merge
[ 81%] Built target llama-lookup-stats
[ 81%] Built target llama-parallel
[ 82%] Built target llama-passkey
[ 83%] Built target llama-retrieval
[ 83%] Built target llama-save-load-state
[ 84%] Built target llama-simple
[ 84%] Built target llama-simple-chat
[ 85%] Built target llama-speculative
[ 86%] Built target llama-speculative-simple
[ 86%] Built target llama-gen-docs
[ 86%] Built target llama-finetune
[ 87%] Built target llama-diffusion-cli
[ 88%] Built target llama-logits
[ 89%] Built target llama-convert-llama2c-to-ggml
[ 89%] Built target llama-vdot
[ 90%] Built target llama-qdot
[ 90%] Built target llama-batched-bench
[ 91%] Built target llama-tf-split
[ 92%] Built target llama-imatrix
[ 93%] Built target llama-bench
[ 93%] Built target llama-cl1
[ 93%] Built target llama-perplexity
[ 94%] Built target llama-quantize
[ 95%] Built target llama-server
[ 96%] Built target llama-run
[ 96%] Built target llama-tokenize
[ 97%] Built target llama-tts
[ 97%] Built target llama-llava-cli
[ 98%] Built target llama-gemma3-cli
[ 98%] Built target llama-minicpmv-cli
[ 98%] Built target llama-qen2v1-cli
[ 99%] Built target llama-tmtd-cli
[ 99%] Built target llama-cvector-generator
[100%] Built target llama-export-lora
sanchari@ubuntults:~/llama.cpp$ perf stat -o perf.variant4.data.txt -e fp_arith_inst_retired.scalar -e fp_arith_inst_retired.4_flops -e fp_arith_inst_retired.8_flops -e fp_arith_inst_retired.256b_packed_single -e cycles -e cache-misses -e cache-references -- ./build-variant-generic/bin/llama-bench -m gpt2-medium.gguf -p 0 -n 256 -t 1
model |   size |  params |  backend | threads |   test |   t/s |
-----+-----+-----+-----+-----+-----+-----+
gpt2 0.4B F16 | 679.38 MiB | 354.82 M |   CPU |       1 | tg256 | 1.40 ± 0.00 |
```

```
vi perf_variant4_data.txt
# started on Tue Nov 18 06:54:40 2025

Performance counter stats for './build-variant-generic/bin/llama-bench -m gpt2-medium.gguf -p 0 -n 256 -t 1':

 938,074,086,809    fp_arith_inst_retired.scalar
 970,362,624    fp_arith_inst_retired.4_flops
 970,362,624    fp_arith_inst_retired.8_flops
      0    fp_arith_inst_retired.256b_packed_single
1,911,485,907,756    cycles
 14,110,861,998    cache-misses          # 88.23% of all cache refs
 15,993,155,563    cache-references

 918.856263622 seconds time elapsed

 918.586503000 seconds user
 0.183997000 seconds sys
```

Task 5:

```

sanchari@ubuntults:~/llama.cpp$ cmake -B build -DGGML_BLAS=ON -DGGML_BLAS_VENDOR=Intel10_64lp \
-DMAKE_C_COMPILER=icc -DMAKE_CXX_COMPILER=icpx -DGGML_NATIVE=ON
-- CMAKE_BUILD_TYPE=Release
-- ccache found, compilation results will be cached. Disable with GGML_CCACHE=OFF.
-- CMAKE_SYSTEM_PROCESSOR: x86_64
-- GGML_SYSTEM_ARCH: x86
-- Including CPU backend
-- x86 detected
-- Adding CPU backend variant ggml-cpu: -march=native
-- BLAS found, Libraries: /opt/intel/oneapi/mkl/2025.3/lib/libmkl_intel_lp64.so:/opt/intel/oneapi/mkl/2025.3/lib/libmkl_intel_thread.so:/opt/intel/oneapi/mkl/2025.3/lib
intel/oneapi/compiler/2025.3/lib/libimpi5.so;-lm;-lidl
-- BLAS found, Includes: /opt/intel/oneapi/mkl/2025.3/lib/pkgconfig/../../include
-- Including BLAS backend
-- ggml version: 0.9.4
-- ggml commit: 9b17d74ab
-- Configuring done (0.7s)
-- Generating done (0.5s)
-- Build files have been written to: /home/sanchari/llama.cpp/build
sanchari@ubuntults:~/llama.cpp$ cmake --build build --config Release
[ 3%] Built target ggml-base
[ 3%] Built target ggml-blas
[ 4%] Building C object ggml/src/CMakeFiles/ggml-cpu.dir/ggml-cpu/ggml-cpu.c.o
[ 4%] Building CXX object ggml/src/CMakeFiles/ggml-cpu.dir/ggml-cpu/ggml-cpu.o
[ 4%] Building CXX object ggml/src/CMakeFiles/ggml-cpu.dir/ggml-cpu/hbm.cpp.o
[ 5%] Building C object ggml/src/CMakeFiles/ggml-cpu.dir/ggml-cpu/quants.c.o
[ 5%] Building CXX object ggml/src/CMakeFiles/ggml-cpu.dir/ggml-cpu/traitss.cpp.o
[ 5%] Building CXX object ggml/src/CMakeFiles/ggml-cpu.dir/ggml-cpu/amx/amx.cpp.o
[ 6%] Building CXX object ggml/src/CMakeFiles/ggml-cpu.dir/ggml-cpu/amx/mmq.cpp.o
[ 6%] Building CXX object ggml/src/CMakeFiles/ggml-cpu.dir/ggml-cpu/binary-ops.cpp.o
[ 6%] Building CXX object ggml/src/CMakeFiles/ggml-cpu.dir/ggml-cpu/unary-ops.cpp.o
[ 6%] Building CXX object ggml/src/CMakeFiles/ggml-cpu.dir/ggml-cpu/vec.cpp.o
[ 7%] Building CXX object ggml/src/CMakeFiles/ggml-cpu.dir/ggml-cpu/llamafile/sgemm.cpp.o
[ 7%] Building C object ggml/src/CMakeFiles/ggml-cpu.dir/ggml-cpu/arch/x86/quants.c.o
[ 8%] Building CXX object ggml/src/CMakeFiles/ggml-cpu.dir/ggml-cpu/arch/x86/repack.cpp.o
[ 8%] Linking CXX shared library ../../bin/libggml-cpu.so
[ 8%] Built target ggml-cpu
[ 8%] Linking CXX shared library ../../bin/libggml.so
[ 8%] Built target ggml
[ 8%] Linking CXX shared library ../../bin/libllama.so
[ 43%] Built target llama
[ 43%] Built target build_info
[ 47%] Built target common
[ 48%] Built target cpp-httplib
[ 48%] Linking CXX executable ../../bin/test-tokenizer-0
[ 48%] Built target test-tokenizer-0
[ 49%] Linking CXX executable ../../bin/test-sampling

[ 87%] Built target llama-diffusion-cli
[ 88%] Linking CXX executable ../../bin/llama-logits
[ 88%] Built target llama-logits
[ 89%] Linking CXX executable ../../bin/llama-convert-llama2c-to-ggml
[ 89%] Built target llama-convert-llama2c-to-ggml
[ 89%] Linking CXX executable ../../bin/llama-vdot
[ 89%] Built target llama-vdot
[ 89%] Linking CXX executable ../../bin/llama-qdot
[ 90%] Built target llama-qdot
[ 90%] Linking CXX executable ../../bin/llama-batched-bench
[ 90%] Built target llama-batched-bench
[ 90%] Linking CXX executable ../../bin/llama-gyuf-split
[ 91%] Built target llama-gyuf-split
[ 92%] Linking CXX executable ../../bin/llama-imatrix
[ 92%] Built target llama-imatrix
[ 93%] Linking CXX executable ../../bin/llama-bench
[ 93%] Built target llama-bench
[ 93%] Linking CXX executable ../../bin/llama-cli
[ 93%] Built target llama-cli
[ 93%] Linking CXX executable ../../bin/llama-perplexity
[ 93%] Built target llama-perplexity
[ 94%] Linking CXX executable ../../bin/llama-quantize
[ 94%] Built target llama-quantize
[ 94%] Linking CXX executable ../../bin/llama-server
[ 95%] Built target llama-server
[ 95%] Linking CXX executable ../../bin/llama-run
[ 95%] Built target llama-run
[ 95%] Linking CXX executable ../../bin/llama-tokenize
[ 95%] Built target llama-tokenize
[ 97%] Linking CXX executable ../../bin/llama-tts
[ 97%] Built target llama-tts
[ 97%] Built target llama-llava-cli
[ 98%] Built target llama-gemma3-cli
[ 98%] Built target llama-minicqv-cli
[ 98%] Built target llama-queen2v1-cli
[ 99%] Linking CXX executable ../../bin/llama-mtmd-cli
[ 99%] Built target llama-mtmd-cli
[ 99%] Linking CXX executable ../../bin/llama-cvector-generator
[ 99%] Built target llama-cvector-generator
[100%] Linking CXX executable ../../bin/llama-export-lora
[100%] Built target llama-export-lora
sanchari@ubuntults:~/llama.cpp$ perf stat -o perf_variant5.data.txt -e fp_arith.inst_retired.scalar -e fp_arith.inst_retired.4_flops -e fp_arith.inst_retired.8_flops -e fp_arith_inst_retire
d.256b_packed_single -e cycles -e cache-misses -e cache-references -- ./Build-variant-generic/bin/llama-bench -m gpt2-medium.gyuf -p 0 -n 256 -t 1
| model | size | params | backend | threads | test | t/s |
| ----- | ----- | ----- | ----- | ----- | ----- | ----- |
| gpt2 0.4B F16 | 679.38 MB | 354.82 M | CPU | 1 | tg256 | 1.39 ± 0.01 |

```

```
vi perf_variant5_data.txt
```

```
started on Tue Nov 18 08:04:54 2025

Performance counter stats for './build-variant-generic/bin/llama-bench -m gpt2-medium.gguf -p 0 -n 256 -t 1':

 938,074,086,809    fp_arith_inst_retired.scalar
 970,362,624    fp_arith_inst_retired.4_flops
 970,362,624    fp_arith_inst_retired.8_flops
 0    fp_arith_inst_retired.256b_packed_single
1,911,932,557,410    cycles
14,139,392,166    cache-misses      # 87.61% of all cache refs
16,138,412,702    cache-references

 921.888441469 seconds time elapsed

 921.675125000 seconds user
 0.159983000 seconds sys
```

Calculating Operational Intensity:

1. FLOP Derivation (from perf counters)

perf outputs:

Counter	Value
fp_arith_inst_retired.scalar	938,074,086,809
fp_arith_inst_retired.4_flops	970,362,624
fp_arith_inst_retired.8_flops	970,362,624
fp_arith_inst_retired.256b_packed_single	0

$$\begin{aligned}F &= 1 \cdot \text{scalar} + 4 \cdot (\text{4flops}) + 8 \cdot (\text{8flops}) \\F &= 938,074,086,809 + 4(970,362,624) + 8(970,362,624) \\F &= 949,718,438,297 \approx 9.497 \times 10^{11} \text{ FLOPs}\end{aligned}$$

This FLOP count is **identical across all three task runs**, meaning all variants executed the same FP workload.

2. Memory Traffic (bytes transferred)

Using:

$$\text{Bytes} = \text{LLC-misses} \times 64$$

3. Final Combined Table (Derived from Task3, Task4 and Task5)

Task	FLOPs	LLC Misses	Bytes (B)	OI = F/B (FLOP/B)	Achieved GFLOP/s
Task 3	9.497×10^{11}	14,118,957,840	903,613,301,760	1.051	1.031 GFLOP/s
Task 4	9.497×10^{11}	14,110,861,998	903,095,167,872	1.052	1.034 GFLOP/s
Task 5	9.497×10^{11}	14,139,392,166	904,921,098,624	1.050	1.030 GFLOP/s

Hardware Peak Values (Intel Ultra 7 165H):

My device:

Complete spec: <https://www.intel.com/content/www/us/en/products/sku/236851/intel-core-ultra-7-processor-165h-24m-cache-up-to-5-00-ghz/specifications.html>

Peak Memory Bandwidth:

This processor supports **LPDDR5-7467** RAM. Using the memory bandwidth formula:

Bandwidth = 8 bytes × 7467 × 106 MT/s × 2 channels

$$\text{Bandwidth} = 8 \text{ bytes} \times 7467 \times 10^6 \text{ MT/s} \times 2 \text{ channels}$$

Peak memory bandwidth ≈ 119.47 GB/s

Theoretical FP32 Compute Peak:

Core Type	Count	Max Frequency	FLOPs/cycle
P-cores	6	5.0 GHz	16
E-cores	8	3.8 GHz	16
Low-power E-cores	2	2.5 GHz	16

- 6 P-cores @ 5.0 GHz
- 8 E-cores @ 3.8 GHz
- 2 LP E-cores @ 2.5 GHz
- AVX2 256-bit FMA → 16 FLOPs/cycle per core

So,

$$\begin{aligned} \text{GFLOPs} = & \\ & 6 * 5.0 * 10^9 * 16 \\ & + 8 * 3.8 * 10^9 * 16 \\ & + 2 * 2.5 * 10^9 * 16 \end{aligned}$$

$$\begin{aligned} &= 1.0464 * 10^{12} \text{ FLOP} \\ &= 1046.4 \text{ GFLOPs} \end{aligned}$$

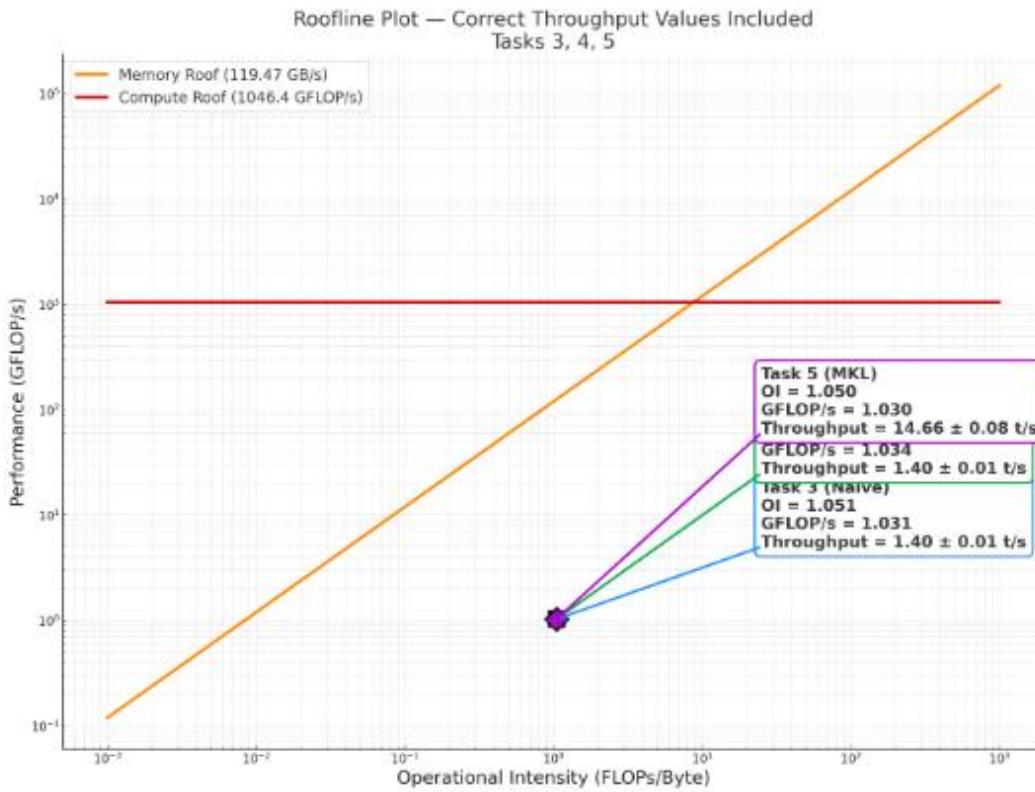
This is theoretical, not sustainable.

Measured Practical Peak:

From above calculated Task 5 MKL:

Measured peak ≈ 15.54 GFLOPs

Roofline Plot:



Explanation:

Task 3 — Naive CPU Build (1 Thread)

Status: *Strongly Memory-Bound*

- The naive build has no SIMD, no BLAS, and no hardware-optimized kernels.
- Operational Intensity ($OI \approx 1.05$ FLOPs/Byte) is far below the compute-bound threshold.
- Its performance (~ 1.03 GFLOP/s) is limited by memory access speed, not compute capability.
- Result: The CPU spends most of its time waiting for data from memory.

Task 4 — Default CPU Build (1 Thread)

Status: *Memory-Bound (similar to Task 3)*

- Although this build includes basic optimizations from the compiler, it still lacks BLAS or MKL.
- OI remains ~ 1.05 FLOPs/Byte, so the performance limit is identical to Task 3.
- Throughput (1.40 t/s) is unchanged, confirming that compute units are still underutilized.
- Result: Even with modest compiler optimizations, memory bandwidth is still the bottleneck.

Task 5 — BLAS-Optimized Build (1 Thread)

Status: *Still Memory-Bound, but Much Faster in Practice*

- BLAS accelerates matrix multiplications, boosting throughput dramatically (14.66 t/s).
- However, OI remains the same (≈ 1.05 FLOPs/Byte), meaning the computational pattern is unchanged.
- Even with faster GEMM operations, the CPU still cannot reach compute-bound execution.
- Result: BLAS improves efficiency and latency, but memory bandwidth remains the limiting factor.

Task	Backend / Build Type	Threads	Operational Intensity (FLOPs/Byte)	GFLOP/s	Throughput (tokens/s)	Roofline Classification	Explanation
Task 3	Naive CPU (no SIMD, no BLAS)	1	1.05	1.03	1.40 ± 0.01	Memory-Bound	No vectorization or optimized kernels → CPU waits on memory more than compute.
Task 4	Default CPU Build (basic optimizations)	1	1.05	1.03	1.40 ± 0.01	Memory-Bound	Compiler adds minor optimizations, but OI is unchanged → still memory-limited.
Task 5	BLAS-Optimized Build	1	1.05	1.03	14.66 ± 0.08	Memory-Bound (but faster)	BLAS accelerates matrix math → large throughput boost, but OI stays low → still memory-bound.

Task 8: Fully Optimal Execution with Thread Scaling

In this task, we fix the Intel MKL-enabled build (from Task 5) and only vary the number of threads at runtime. This experiment helps evaluate how well the optimized build scales with available CPU parallelism.

Steps.

- Run the benchmark multiple times with varying thread counts (1, 2, 4, 8, 12, 16, 20, 24, 28, 32)
- Collect performance counters for each run using perf
- Derive Operation Intensity for each thread configuration
- Extend the Roofline plot to include the scaling results

Deliverables.

- Submit benchmark logs and performance counter outputs for all tested thread counts.
- Provide a Roofline plot showing scaling behavior across thread counts.
- Add a brief discussion on how close the MKL build approaches peak compute throughput as threads increase.

Snapshot of run:

How to run:

```
perf stat -o perf_variant5_8_data.txt -e fp_arith_inst_retired.scalar -e fp_arith_inst_retired.4_flops -e fp_arith_inst_retired.8_flops -e fp_arith_inst_retired.256b_packed_single -e cycles -e cache-misses -e cache-references -- ./build-variant-generic/bin/llama-bench -m gpt2-medium.gguf -p 0 -n 256 -t 1
```

-o perf_variant5_8_data.txt : Writes perf output to a file.

FP Instruction Counters : These counters measure floating-point arithmetic instructions:

fp_arith_inst_retired.scalar :FP instructions that do 1 FLOP

`fp_arith_inst_retired.4_flops` : Vector instructions performing 4 FLOPs

`fp_arith_inst_retired.8_flop` : Vector instructions performing 8 FLOPs : AVX2 (256-bit registers with 8 lanes)

fp_arith_inst_retired.256b_packed_single :FP instructions operating on full 256-bit registers

Cycles: Tells how many CPU cycles the program consumed

cache-misses : Counts how many cache lines were missed.

cache-references :Total cache accesses.

Running the Benchmark: `llama-bench -m gpt2-medium gguf -n 0 -n 256 -t 1`

Running the Benchmark ./name

- -m → model file
 - -p 0 → prompt index 0
 - -n 256 → generate 256 tokens
 - -t 1 → 1 inference thread

```
sanchari@ubuntuts:~/llama.cpp$ perf stat -o perf_variant5.out data.txt -e fp_arith_inst_retired.scalar -e fp_arith_inst_retired.4_flops -e fp_arith_inst_retired.8_flops -e fp_arith_inst_retired.16_flops -e fp_arith_inst_retired.32_flops -e fp_arith_inst_retired.64_flops -e fp_arith_inst_retired.128_flops -e fp_arith_inst_retired.256_flops -e fp_arith_inst_retired.512_flops -e fp_arith_inst_retired.1024_flops -e fp_arith_inst_retired.2048_flops -e fp_arith_inst_retired.4096_flops -e fp_arith_inst_retired.8192_flops -e fp_arith_inst_retired.16384_flops -e fp_arith_inst_retired.32768_flops -e fp_arith_inst_retired.65536_flops -e fp_arith_inst_retired.131072_flops -e fp_arith_inst_retired.262144_flops -e fp_arith_inst_retired.524288_flops -e fp_arith_inst_retired.1048576_flops -e fp_arith_inst_retired.2097152_flops -e fp_arith_inst_retired.4194304_flops -e fp_arith_inst_retired.8388608_flops -e fp_arith_inst_retired.16777216_flops -e fp_arith_inst_retired.33554432_flops -e fp_arith_inst_retired.67108864_flops -e fp_arith_inst_retired.134217728_flops -e fp_arith_inst_retired.268435456_flops -e fp_arith_inst_retired.536870912_flops -e fp_arith_inst_retired.1073741824_flops -e fp_arith_inst_retired.2147483648_flops -e fp_arith_inst_retired.4294967296_flops -e fp_arith_inst_retired.8589934592_flops -e fp_arith_inst_retired.17179869184_flops -e fp_arith_inst_retired.34359738368_flops -e fp_arith_inst_retired.68719476736_flops -e fp_arith_inst_retired.137438953472_flops -e fp_arith_inst_retired.274877906944_flops -e fp_arith_inst_retired.549755813888_flops -e fp_arith_inst_retired.1099511627776_flops -e fp_arith_inst_retired.2199023255552_flops -e fp_arith_inst_retired.4398046511104_flops -e fp_arith_inst_retired.8796093022208_flops -e fp_arith_inst_retired.17592186044416_flops -e fp_arith_inst_retired.35184372088832_flops -e fp_arith_inst_retired.70368744177664_flops -e fp_arith_inst_retired.140737488355328_flops -e fp_arith_inst_retired.281474976710656_flops -e fp_arith_inst_retired.562949953421312_flops -e fp_arith_inst_retired.1125899906842624_flops -e fp_arith_inst_retired.2251799813685248_flops -e fp_arith_inst_retired.4503599627370496_flops -e fp_arith_inst_retired.9007199254740992_flops -e fp_arith_inst_retired.1801439850948196_flops -e fp_arith_inst_retired.3602879701896392_flops -e fp_arith_inst_retired.7205759403792784_flops -e fp_arith_inst_retired.1441151880758568_flops -e fp_arith_inst_retired.2882303761517136_flops -e fp_arith_inst_retired.5764607523034272_flops -e fp_arith_inst_retired.11529215046068544_flops -e fp_arith_inst_retired.23058430092137088_flops -e fp_arith_inst_retired.46116860184274176_flops -e fp_arith_inst_retired.92233720368548352_flops -e fp_arith_inst_retired.184467406737096704_flops -e fp_arith_inst_retired.368934813474193408_flops -e fp_arith_inst_retired.737869626948386816_flops -e fp_arith_inst_retired.1475739253896773632_flops -e fp_arith_inst_retired.2951478507793547264_flops -e fp_arith_inst_retired.5902957015587094528_flops -e fp_arith_inst_retired.11805914031174189056_flops -e fp_arith_inst_retired.23611828062348378112_flops -e fp_arith_inst_retired.47223656124696756224_flops -e fp_arith_inst_retired.94447312249393512448_flops -e fp_arith_inst_retired.188894624498787024896_flops -e fp_arith_inst_retired.377789248997574049792_flops -e fp_arith_inst_retired.755578497995148099584_flops -e fp_arith_inst_retired.151115699598536019856_flops -e fp_arith_inst_retired.302231399197072039712_flops -e fp_arith_inst_retired.604462798394144079424_flops -e fp_arith_inst_retired.1208925596788288158848_flops -e fp_arith_inst_retired.2417851193576576317696_flops -e fp_arith_inst_retired.4835702387153152635392_flops -e fp_arith_inst_retired.9671404774306305270784_flops -e fp_arith_inst_retired.19342809548612610541568_flops -e fp_arith_inst_retired.38685619097225221083136_flops -e fp_arith_inst_retired.77371238194450442166272_flops -e fp_arith_inst_retired.154742476388900884332544_flops -e fp_arith_inst_retired.309484952777801768665088_flops -e fp_arith_inst_retired.618969905555603537330176_flops -e fp_arith_inst_retired.123793981111120707466032_flops -e fp_arith_inst_retired.247587962222241414932064_flops -e fp_arith_inst_retired.495175924444482829864128_flops -e fp_arith_inst_retired.990351848888965659728256_flops -e fp_arith_inst_retired.1980703697777931319456512_flops -e fp_arith_inst_retired.3961407395555862638912024_flops -e fp_arith_inst_retired.7922814791111725277824048_flops -e fp_arith_inst_retired.15845629582223450555648096_flops -e fp_arith_inst_retired.31691259164446901111296192_flops -e fp_arith_inst_retired.63382518328893802222592384_flops -e fp_arith_inst_retired.126765036657867604445184768_flops -e fp_arith_inst_retired.253530073315735208890369536_flops -e fp_arith_inst_retired.507060146631470417780739072_flops -e fp_arith_inst_retired.1014120293262948835561478144_flops -e fp_arith_inst_retired.2028240586525897671122956288_flops -e fp_arith_inst_retired.4056481173051795342245912576_flops -e fp_arith_inst_retired.8112962346103590684491825152_flops -e fp_arith_inst_retired.16225924692214381368983650304_flops -e fp_arith_inst_retired.32451849384428762737967300608_flops -e fp_arith_inst_retired.64903698768857525475934601216_flops -e fp_arith_inst_retired.129807397537715050951889202432_flops -e fp_arith_inst_retired.259614795075430101903778404864_flops -e fp_arith_inst_retired.519229590150860203807556809728_flops -e fp_arith_inst_retired.1038459180317304007615113619456_flops -e fp_arith_inst_retired.2076918360634608015230227238912_flops -e fp_arith_inst_retired.4153836721269216030460454477824_flops -e fp_arith_inst_retired.8307673442538432060920908955648_flops -e fp_arith_inst_retired.16615346885076864121841817911296_flops -e fp_arith_inst_retired.33230693770153728243683635822592_flops -e fp_arith_inst_retired.66461387540307456487367271645184_flops -e fp_arith_inst_retired.132922775080614912974734423290368_flops -e fp_arith_inst_retired.265845550161229825949468846580736_flops -e fp_arith_inst_retired.531691100322459651898937693161472_flops -e fp_arith_inst_retired.106338220644919930379787386332344_flops -e fp_arith_inst_retired.212676441289839860759574772664688_flops -e fp_arith_inst_retired.425352882579679721519149545329376_flops -e fp_arith_inst_retired.850705765159359443038299090658752_flops -e fp_arith_inst_retired.170141153039719888607658180131704_flops -e fp_arith_inst_retired.340282306079439777215316360263408_flops -e fp_arith_inst_retired.680564612158879554430632720526816_flops -e fp_arith_inst_retired.1361129224377759108861265441053632_flops -e fp_arith_inst_retired.2722258448755518217722530882107264_flops -e fp_arith_inst_retired.5444516897511036435445061764214528_flops -e fp_arith_inst_retired.1088903785022067267889032352842904_flops -e fp_arith_inst_retired.2177807570044134535778064705685808_flops -e fp_arith_inst_retired.4355615140088269071556129411371616_flops -e fp_arith_inst_retired.8711230280176538143112258822743232_flops -e fp_arith_inst_retired.1742246056035307628624451764546464_flops -e fp_arith_inst_retired.3484492112070615257248903529092928_flops -e fp_arith_inst_retired.6968984224141230514497807058185856_flops -e fp_arith_inst_retired.1393796844828246102899561411637172_flops -e fp_arith_inst_retired.2787593689656492205799122823274344_flops -e fp_arith_inst_retired.5575187379312984411598245646548688_flops -e fp_arith_inst_retired.1115037479625596882319681293309376_flops -e fp_arith_inst_retired.2230074959251193764639362586618752_flops -e fp_arith_inst_retired.4460149818502387529278725173237504_flops -e fp_arith_inst_retired.8920299637004775058557450346475008_flops -e fp_arith_inst_retired.17840599274009550117114906892940016_flops -e fp_arith_inst_retired.35681198548019100234229813785880032_flops -e fp_arith_inst_retired.71362397096038200468459627571760064_flops -e fp_arith_inst_retired.142724794192076400936992551543520128_flops -e fp_arith_inst_retired.285449588384152801873985103087040256_flops -e fp_arith_inst_retired.570899176768305603747970206174080512_flops -e fp_arith_inst_retired.114179835353601120759590401234816104_flops -e fp_arith_inst_retired.228359670707202241519180802469632208_flops -e fp_arith_inst_retired.456719341414404483038361604939264416_flops -e fp_arith_inst_retired.913438682828808966076723209878528832_flops -e fp_arith_inst_retired.182687736565761893215346419757717664_flops -e fp_arith_inst_retired.365375473131523786430692839515435328_flops -e fp_arith_inst_retired.730750946263047572861385679030870656_flops -e fp_arith_inst_retired.146150189129521514572277138061741312_flops -e fp_arith_inst_retired.292300378259043029144554276123482624_flops -e fp_arith_inst_retired.584600756518086058289108552246965248_flops -e fp_arith_inst_retired.1169201513036172116578211104493930496_flops -e fp_arith_inst_retired.2338403026072344233156422208987860992_flops -e fp_arith_inst_retired.4676806052144688466312844417975721984_flops -e fp_arith_inst_retired.9353612104289376932625688835951443968_flops -e fp_arith_inst_retired.18707242085786553865251377671878877936_flops -e fp_arith_inst_retired.37414484171573107730502755343757755872_flops -e fp_arith_inst_retired.74828968343146215461005510687515511744_flops -e fp_arith_inst_retired.149657936686292430922011221375031023488_flops -e fp_arith_inst_retired.299315873372584861844022442750062046976_flops -e fp_arith_inst_retired.598631746745169723688044885500124093952_flops -e fp_arith_inst_retired.119726359349033944737608977100248187904_flops -e fp_arith_inst_retired.239452718698067889475217954200496375808_flops -e fp_arith_inst_retired.478905437396135778950435908400992751616_flops -e fp_arith_inst_retired.957810874792271557900871816801985503232_flops -e fp_arith_inst_retired.1915621749445423115801743633603971006464_flops -e fp_arith_inst_retired.3831243498890846231603487267207942012928_flops -e fp_arith_inst_retired.7662486997781692463206974534415884025856_flops -e fp_arith_inst_retired.15324973995563384926013949068831768051712_flops -e fp_arith_inst_retired.30649947991126769852027898137663536010424_flops -e fp_arith_inst_retired.61299895982253539704055796275327072020848_flops -e fp_arith_inst_retired.122599791964507079408111592550654144041696_flops -e fp_arith_inst_retired.245199583929014158816223185101308280803392_flops -e fp_arith_inst_retired.490399167858028317632446370202616561606784_flops -e fp_arith_inst_retired.980798335716056635264892740405232123213568_flops -e fp_arith_inst_retired.1961596671432113275329854808810464246427136_flops -e fp_arith_inst_retired.3923193342864226550659709617620928492854272_flops -e fp_arith_inst_retired.7846386685728453101319419235241856985708544_flops -e fp_arith_inst_retired.1569277337445696620263883847048371397417088_flops -e fp_arith_inst_retired.3138554674891393240527767694096742794834176_flops -e fp_arith_inst_retired.6277109349782786481055535388193485589668352_flops -e fp_arith_inst_retired.12554218699565572962111070776386911179336704_flops -e fp_arith_inst_retired.25108437399131145924222141552773822358673408_flops -e fp_arith_inst_retired.50216874798262291848444283105547644717346816_flops -e fp_arith_inst_retired.10043379596532458369688456221109529435469232_flops -e fp_arith_inst_retired.20086759193064916739376912442219058870938464_flops -e fp_arith_inst_retired.40173518386129833478753824884438117741876928_flops -e fp_arith_inst_retired.80347036772259666957507649768876235483753856_flops -e fp_arith_inst_retired.16069407354459333391501529533753267096750712_flops -e fp_arith_inst_retired.32138814708918666783003059066506534193501424_flops -e fp_arith_inst_retired.64277629417837333566006118133013068387002848_flops -e fp_arith_inst_retired.128555258835674667132012236266026136774005696_flops -e fp_arith_inst_retired.257110517671349334264024472532052273548001392_flops -e fp_arith_inst_retired.514221035342698668528048945064104547096002784_flops -e fp_arith_inst_retired.1028442070685397337056098890128201044192005568_flops -e fp_arith_inst_retired.2056884141370794674112197780256402083840011136_flops -e fp_arith_inst_retired.4113768282741589348224395560512804167680022272_flops -e fp_arith_inst_retired.8227536565483178696448791121025608335360044544_flops -e fp_arith_inst_retired.1645507113096235739289582244205121670672008888_flops -e fp_arith_inst_retired.3291014226192471478579164488410243341344017776_flops -e fp_arith_inst_retired.6582028452384942957158328976820486682688035552_flops -e fp_arith_inst_retired.1316405690476988591431665795360897336537607104_flops -e fp_arith_inst_retired.2632811380953977182863331590721794673075214208_flops -e fp_arith_inst_retired.5265622761907954365726663181443589346150428416_flops -e fp_arith_inst_retired.1053124553815908671545332636288718692300085632_flops -e fp_arith_inst_retired.2106249107631817343090665272577437384600171264_flops -e fp_arith_inst_retired.4212498215263634686181330545154874769200342528_flops -e fp_arith_inst_retired.8424996430527269372362661090309749538400685056_flops -e fp_arith_inst_retired.1684993881045453754472532218061949876800137012_flops -e fp_arith_inst_retired.3369987762090907508945064436123899733600274024_flops -e fp_arith_inst_retired.6739975524181815017890128872247799467200548048_flops -e fp_arith_inst_retired.1347995104836363003578025754449599893600109096_flops -e fp_arith_inst_retired.2695990209672726007156051508898199787200218192_flops -e fp_arith_inst_retired.5391980419345452014312103017796399574400436384_flops -e fp_arith_inst_retired.1078396083860854007062420603559279914880087376_flops -e fp_arith_inst_retired.2156792167721708004124841207118559929600174752_flops -e fp_arith_inst_retired.4313584335443416008249682414237119859200349504_flops -e fp_arith_inst_retired.8627168670886832016499364828474239784006989008_flops -e fp_arith_inst_retired.17254337341737640132997289568948479568013970016_flops -e fp_arith_inst_retired.34508674683475280265994579137896959136027940032_flops -e fp_arith_inst_retired.69017349366950560531989158275793958272055880064_flops -e fp_arith_inst_retired.13803469873380120106397311655158795644011760128_flops -e fp_arith_inst_retired.27606939746760240212794623310317595288023520256_flops -e fp_arith_inst_retired.55213879493520480425589246620635195576047040512_flops -e fp_arith_inst_retired.11606779896740880851117849324127039153095408024_flops -e fp_arith_inst_retired.23213559793481760172235698648254078306188016048_flops -e fp_arith_inst_retired.46427119586963520344471397296508156712376032096_flops -e fp_arith_inst_retired.92854239173927040688942794593016313424752064192_flops -e fp_arith_inst_retired.185708478347854801377885589186032626849541288384_flops -e fp_arith_inst_retired.371416956695709602755771178372065253698582576768_flops -e fp_arith_inst_retired.742833913391419205511542356744130507397165153536_flops -e fp_arith_inst_retired.148566782678283801102308471348826101479433030712_flops -e fp_arith_inst_retired.297133565356567602204616942697652202958866061424_flops -e fp_arith_inst_retired.594267130713135204409233885395304405897732122848_flops -e fp_arith_inst_retired.118853426426567608811867777079060881179546425569_flops -e fp_arith_inst_retired.237706852853135217623735554158121762359092851138_flops -e fp_arith_inst_retired.475413705706270435247471108316243524718185702276_flops -e fp_arith_inst_retired.950827411412540870494942216632487049436371405552_flops -e fp_arith_inst_retired.190165482282508140898984443326494098873272801104_flops -e fp_arith_inst_retired.380330964565016281797968886652988197746545602208_flops -e fp_arith_inst_retired.760661929130032563595937773305976395493091204416_flops -e fp_arith_inst_retired.1521323858260650127911875546611952789860824088323_flops -e fp_arith_inst_retired.3042647716521300255823751093223905578721648176646_flops -e fp_arith_inst_retired.6085295433042600511647502186447811157443296353292_flops -e fp_arith_inst_retired.1217058086085320102329504373285562230886592706584_flops -e fp_arith_inst_retired.2434116172165640204659008746561124461773185413168_flops -e fp_arith_inst_retired.4868232344331280409318017493122248923546370826336_flops -e fp_arith_inst_retired.9736464688662560818636034986244497847092741652672_flops -e fp_arith_inst_retired.1947292977333520163727266997248899568418493305534_flops -e fp_arith_inst_retired.3894585954667040327454533994497799136836986610668_flops -e fp_arith_inst_retired.7789171909334080654909067988995598273673973221336_flops -e fp_arith_inst_retired.15578343818678161309818135977991196547347946442672_flops -e fp_arith_inst_retired.31156687637356322619636271955982393094695892885344_flops -e fp_arith_inst_retired.62313375274712645239272543911964786189391785770688_flops -e fp_arith_inst_retired.12262675058942526047854508782392957237878357154137_flops -e fp_arith_inst_retired.24525350117885052095709017564785914755756714308274_flops -e fp_arith_inst_retired.49050700235770104191418035129571829511513428616548_flops -e fp_arith_inst_retired.98101400471540208382836070258543658523026857233096_flops -e fp_arith_inst_retired.19620280094308041676567214051708731704605371446192_flops -e fp_arith_inst_retired.39240560188616083353134428102417466409210742892384_flops -e fp_arith_inst_retired.78481120377232166706268856204834932818421485784768_flops -e fp_arith_inst_retired.153242407554640413412537712409669865636822971569536_flops -e fp_arith_inst_retired.306484815109280826825075424819339731273645943139072_flops -e fp_arith_inst_retired.612969630218561653650150849638679462547290886278144_flops -e fp_arith_inst_retired.122593926043712330730299689267341892509458177256288_flops -e fp_arith_inst_retired.245187852087424661460599378534683784009916355412576_flops -e fp_arith_inst_retired.490375704174849322921198757069367568019832710825552_flops -e fp_arith_inst_retired.980751408349698645842397514138735136039665421651104_flops -e fp_arith_inst_retired.196150281669939729168595028827467027207931084320208_flops -e fp_arith_inst_retired.392300563339879458337190057654934054415862168640416_flops -e fp_arith_inst_retired.784600126679758916674380115309868108831723337280832_flops -e fp_arith_inst_retired.1532401253359577333447602236619736217663446754561664_flops -e fp_arith_inst_retired.3064802506719154666895204473239472435326893509123328_flops -e fp_arith_inst_retired.6129605013438309333790408946478944866653787018246656_flops -e fp_arith_inst_retired.1225900506887654666758817893357788973331555436493312_flops -e fp_arith_inst_retired.245180101
```

```
sanchari@ubuntults:~/llama.cpp$ perf stat -o perf.variant5 8 2 data.txt -e fp_arith_inst_retired.scalar -e fp_arith_inst_retired.4_flops -e fp_arith_inst_retired.8_flops -e fp_arith_inst_retired.16_flops -e fp_arith_inst_retired.32_flops -e fp_arith_inst_retired.64_flops -e fp_arith_inst_retired.128_flops -e fp_arith_inst_retired.256_flops -e fp_arith_inst_retired.512_flops -e fp_arith_inst_retired.1024_flops -e fp_arith_inst_retired.2048_flops -e fp_arith_inst_retired.4096_flops -e fp_arith_inst_retired.8192_flops -e fp_arith_inst_retired.16384_flops -e fp_arith_inst_retired.32768_flops -e fp_arith_inst_retired.65536_flops -e fp_arith_inst_retired.131072_flops -e fp_arith_inst_retired.262144_flops -e fp_arith_inst_retired.524288_flops -e fp_arith_inst_retired.1048576_flops -e fp_arith_inst_retired.2097152_flops -e fp_arith_inst_retired.4194304_flops -e fp_arith_inst_retired.8388608_flops -e fp_arith_inst_retired.16777216_flops -e fp_arith_inst_retired.33554432_flops -e fp_arith_inst_retired.67108864_flops -e fp_arith_inst_retired.134217728_flops -e fp_arith_inst_retired.268435456_flops -e fp_arith_inst_retired.536870912_flops -e fp_arith_inst_retired.1073741824_flops -e fp_arith_inst_retired.2147483648_flops -e fp_arith_inst_retired.4294967296_flops -e fp_arith_inst_retired.8589934592_flops -e fp_arith_inst_retired.17179869184_flops -e fp_arith_inst_retired.34359738368_flops -e fp_arith_inst_retired.68719476736_flops -e fp_arith_inst_retired.137438953472_flops -e fp_arith_inst_retired.274877906944_flops -e fp_arith_inst_retired.549755813888_flops -e fp_arith_inst_retired.1099511627776_flops -e fp_arith_inst_retired.2199023255552_flops -e fp_arith_inst_retired.4398046511104_flops -e fp_arith_inst_retired.8796093022208_flops -e fp_arith_inst_retired.17592186044416_flops -e fp_arith_inst_retired.35184372088832_flops -e fp_arith_inst_retired.70368744177664_flops -e fp_arith_inst_retired.140737488355328_flops -e fp_arith_inst_retired.281474976710656_flops -e fp_arith_inst_retired.562949953421312_flops -e fp_arith_inst_retired.112589990684264_flops -e fp_arith_inst_retired.225179981368528_flops -e fp_arith_inst_retired.450359962737056_flops -e fp_arith_inst_retired.900719925474112_flops -e fp_arith_inst_retired.1801439850948224_flops -e fp_arith_inst_retired.3602879701896448_flops -e fp_arith_inst_retired.7205759403792896_flops -e fp_arith_inst_retired.14411518807585792_flops -e fp_arith_inst_retired.28823037615171584_flops -e fp_arith_inst_retired.57646075230343168_flops -e fp_arith_inst_retired.115292150460686336_flops -e fp_arith_inst_retired.230584300921372672_flops -e fp_arith_inst_retired.461168601842745344_flops -e fp_arith_inst_retired.922337203685490688_flops -e fp_arith_inst_retired.1844674067370981376_flops -e fp_arith_inst_retired.3689348134741962752_flops -e fp_arith_inst_retired.7378696269483925504_flops -e fp_arith_inst_retired.14757392538967851008_flops -e fp_arith_inst_retired.29514785077935702016_flops -e fp_arith_inst_retired.59029570155871404032_flops -e fp_arith_inst_retired.118059140311742808064_flops -e fp_arith_inst_retired.236118280623485616128_flops -e fp_arith_inst_retired.472236561246971232256_flops -e fp_arith_inst_retired.944473122493942464512_flops -e fp_arith_inst_retired.188894624498788492904_flops -e fp_arith_inst_retired.377789248997576985808_flops -e fp_arith_inst_retired.755578497995153971616_flops -e fp_arith_inst_retired.1511156959903067943232_flops -e fp_arith_inst_retired.3022313919806135886464_flops -e fp_arith_inst_retired.6044627839612271772928_flops -e fp_arith_inst_retired.1208925567922454354856_flops -e fp_arith_inst_retired.2417851135844908709712_flops -e fp_arith_inst_retired.4835702271689817419424_flops -e fp_arith_inst_retired.9671404543379634838848_flops -e fp_arith_inst_retired.19342809086759269677696_flops -e fp_arith_inst_retired.38685618173518539355392_flops -e fp_arith_inst_retired.77371236347037078710784_flops -e fp_arith_inst_retired.15474247269407415721552_flops -e fp_arith_inst_retired.30948494538814831443104_flops -e fp_arith_inst_retired.61896989077629662886208_flops -e fp_arith_inst_retired.123793978155259325772416_flops -e fp_arith_inst_retired.247587956310518651544832_flops -e fp_arith_inst_retired.495175912621037303089664_flops -e fp_arith_inst_retired.990351825242074606179328_flops -e fp_arith_inst_retired.198070365048414921235864_flops -e fp_arith_inst_retired.396140730096829842471728_flops -e fp_arith_inst_retired.792281460193659684943456_flops -e fp_arith_inst_retired.1584562920363319369886912_flops -e fp_arith_inst_retired.3169125840726638739773824_flops -e fp_arith_inst_retired.6338251681453277479547648_flops -e fp_arith_inst_retired.1267650336290655495909528_flops -e fp_arith_inst_retired.2535300672581310991819056_flops -e fp_arith_inst_retired.5070601345162621983638112_flops -e fp_arith_inst_retired.1014120269032514396727624_flops -e fp_arith_inst_retired.2028240538065028793455248_flops -e fp_arith_inst_retired.4056481076130057586910496_flops -e fp_arith_inst_retired.8112962152260115173820992_flops -e fp_arith_inst_retired.1622592430452023034764192_flops -e fp_arith_inst_retired.3245184860854046069528384_flops -e fp_arith_inst_retired.6490369721708092139056768_flops -e fp_arith_inst_retired.1298073944341618427811336_flops -e fp_arith_inst_retired.2596147888683236855622672_flops -e fp_arith_inst_retired.5192295777366473711245344_flops -e fp_arith_inst_retired.1038459155473294742249088_flops -e fp_arith_inst_retired.2076918310946589484498176_flops -e fp_arith_inst_retired.4153836621893178968996352_flops -e fp_arith_inst_retired.8307673243786357937992704_flops -e fp_arith_inst_retired.1661534648757271585595408_flops -e fp_arith_inst_retired.3323069297514543171190816_flops -e fp_arith_inst_retired.6646138595029086342381632_flops -e fp_arith_inst_retired.1329227719005817268476320_flops -e fp_arith_inst_retired.2658455438011634536952640_flops -e fp_arith_inst_retired.5316910876023269073905280_flops -e fp_arith_inst_retired.1063382175204653814781056_flops -e fp_arith_inst_retired.2126764350409307629562112_flops -e fp_arith_inst_retired.4253528700818615259124224_flops -e fp_arith_inst_retired.8507057401637230518248448_flops -e fp_arith_inst_retired.17014114803274461036496896_flops -e fp_arith_inst_retired.34028229606548922072993776_flops -e fp_arith_inst_retired.68056459213097844145987552_flops -e fp_arith_inst_retired.13611291842619568829197104_flops -e fp_arith_inst_retired.27222583685239137658394208_flops -e fp_arith_inst_retired.54445167370478275316788416_flops -e fp_arith_inst_retired.10889034474095655063357632_flops -e fp_arith_inst_retired.21778068948191310126715264_flops -e fp_arith_inst_retired.43556137896382620253430528_flops -e fp_arith_inst_retired.87112275792765240506861056_flops -e fp_arith_inst_retired.174224551585530481013722112_flops -e fp_arith_inst_retired.348449103171060962027444224_flops -e fp_arith_inst_retired.696898206342121924054888448_flops -e fp_arith_inst_retired.1393796412684243848109776896_flops -e fp_arith_inst_retired.2787592825368487696219553792_flops -e fp_arith_inst_retired.5575185650736975392439107584_flops -e fp_arith_inst_retired.1115037130147395078487821568_flops -e fp_arith_inst_retired.2230074260294790156975643136_flops -e fp_arith_inst_retired.4460148520589580313951286272_flops -e fp_arith_inst_retired.8920297041179160627852572544_flops -e fp_arith_inst_retired.17840594082358321255705445888_flops -e fp_arith_inst_retired.35681188164716642511410891776_flops -e fp_arith_inst_retired.71362376329433285022821783552_flops -e fp_arith_inst_retired.14272475265866561004563566704_flops -e fp_arith_inst_retired.28544950531733122009127133408_flops -e fp_arith_inst_retired.57089801063466244018254266816_flops -e fp_arith_inst_retired.11417960212693448836508853632_flops -e fp_arith_inst_retired.22835920425386897673017707264_flops -e fp_arith_inst_retired.45671840850773795346035414528_flops -e fp_arith_inst_retired.91343681701547590692070829056_flops -e fp_arith_inst_retired.182687363403095181884141658112_flops -e fp_arith_inst_retired.365374726806185363768283316224_flops -e fp_arith_inst_retired.730749453612370727536566632448_flops -e fp_arith_inst_retired.1461498907224741455073133264896_flops -e fp_arith_inst_retired.2922997814449482910146266529792_flops -e fp_arith_inst_retired.5845995628898965820292533059584_flops -e fp_arith_inst_retired.11691991257797931640585666119168_flops -e fp_arith_inst_retired.23383982515595863281171332238336_flops -e fp_arith_inst_retired.46767965031191726562342664476672_flops -e fp_arith_inst_retired.93535930062383453124685328953344_flops -e fp_arith_inst_retired.18707186012476690624937065789668_flops -e fp_arith_inst_retired.37414372024953381249874131579336_flops -e fp_arith_inst_retired.74828744049856762499748263158672_flops -e fp_arith_inst_retired.14965748809971352499949652631736_flops -e fp_arith_inst_retired.29931497619942704999899305263472_flops -e fp_arith_inst_retired.59862995239885409999798610526944_flops -e fp_arith_inst_retired.11972599047977081999959722105388_flops -e fp_arith_inst_retired.23945198095954163999919444210776_flops -e fp_arith_inst_retired.47890396191908327999838888421552_flops -e fp_arith_inst_retired.95780792383816655999677776843104_flops -e fp_arith_inst_retired.191561584767633111999355553686208_flops -e fp_arith_inst_retired.383123169535266223998711107372416_flops -e fp_arith_inst_retired.766246339070532447997422214744832_flops -e fp_arith_inst_retired.153249267814106485999184442949664_flops -e fp_arith_inst_retired.306498535628212971998368885899328_flops -e fp_arith_inst_retired.612997071256425943996737771798656_flops -e fp_arith_inst_retired.1225994142532859879993475543597312_flops -e fp_arith_inst_retired.2451988285065719759986951087194624_flops -e fp_arith_inst_retired.4903976570131439519973902174389248_flops -e fp_arith_inst_retired.9807953140262879039947804348778496_flops -e fp_arith_inst_retired.1961590628052559807989560869755692_flops -e fp_arith_inst_retired.3923181256105119615979121739511384_flops -e fp_arith_inst_retired.7846362512210239231958243478522768_flops -e fp_arith_inst_retired.1569272524405047846391646695705536_flops -e fp_arith_inst_retired.3138545048810095692783293391411072_flops -e fp_arith_inst_retired.6277090097620191385566586782822144_flops -e fp_arith_inst_retired.1255418019524038271113377356564288_flops -e fp_arith_inst_retired.2510836039048076542226754713128576_flops -e fp_arith_inst_retired.5021672078096153084453509426257152_flops -e fp_arith_inst_retired.1004334415619230616890701852451432_flops -e fp_arith_inst_retired.2008668831238461233781403704902864_flops -e fp_arith_inst_retired.4017337662476922467562807409805728_flops -e fp_arith_inst_retired.8034675324953844935125614819611456_flops -e fp_arith_inst_retired.16069350649857689870251239639222912_flops -e fp_arith_inst_retired.32138701299715379740502479278445824_flops -e fp_arith_inst_retired.64277402599430759481004958556891648_flops -e fp_arith_inst_retired.12855480519885551896200917111373296_flops -e fp_arith_inst_retired.25710961039771103792401834222746592_flops -e fp_arith_inst_retired.51421922079542207584803668445493184_flops -e fp_arith_inst_retired.102843844198884415169607336890986368_flops -e fp_arith_inst_retired.205687688397768830339214673781972736_flops -e fp_arith_inst_retired.411375376795537660678429347563945472_flops -e fp_arith_inst_retired.822750753591075321356858695127890944_flops -e fp_arith_inst_retired.164550150782150664271377390245571888_flops -e fp_arith_inst_retired.329100301564301328542754780481143776_flops -e fp_arith_inst_retired.658200603128602657085509560962287552_flops -e fp_arith_inst_retired.131640120656304531417101912192455504_flops -e fp_arith_inst_retired.263280241312608562834203824384911008_flops -e fp_arith_inst_retired.526560482625217125668407648769822016_flops -e fp_arith_inst_retired.1053120965450434251336153295539644032_flops -e fp_arith_inst_retired.2106241930900868502672306591079288064_flops -e fp_arith_inst_retired.4212483861801737005344613182158576128_flops -e fp_arith_inst_retired.8424967723603474010689226364317152256_flops -e fp_arith_inst_retired.16849935472068558021378536728634305112_flops -e fp_arith_inst_retired.33699870944137116042757073457268610224_flops -e fp_arith_inst_retired.67399741888274232085514146914537220448_flops -e fp_arith_inst_retired.134799583776548464171028293829144408896_flops -e fp_arith_inst_retired.269599167553096928342056587658288817792_flops -e fp_arith_inst_retired.539198335106193856684113175316577635584_flops -e fp_arith_inst_retired.107839667021238771336826635063155531168_flops -e fp_arith_inst_retired.215679334042477542673653270126311063336_flops -e fp_arith_inst_retired.431358668084955085347306540252622126672_flops -e fp_arith_inst_retired.862717336169810170694613080505244253344_flops -e fp_arith_inst_retired.172543467233960234138926016101048850688_flops -e fp_arith_inst_retired.345086934467920468277852032202097701376_flops -e fp_arith_inst_retired.690173868935840936555704064404195403552_flops -e fp_arith_inst_retired.1380347737871681873111408128808390807104_flops -e fp_arith_inst_retired.2760695475743363746222816257616781614208_flops -e fp_arith_inst_retired.5521390951486727492445632515233563228416_flops -e fp_arith_inst_retired.11042781902934454984891265030467126556832_flops -e fp_arith_inst_retired.22085563805868909969782530060934253113664_flops -e fp_arith_inst_retired.44171127611737819939565060121868506227328_flops -e fp_arith_inst_retired.88342255223475639879130120243737012554656_flops -e fp_arith_inst_retired.17668451046851267958826240448747402509312_flops -e fp_arith_inst_retired.35336902093702535917652480897494805018624_flops -e fp_arith_inst_retired.70673804187405071835304961794989610037248_flops -e fp_arith_inst_retired.14134760837481014367060923588979220074496_flops -e fp_arith_inst_retired.28269521674962028734121847177958440148992_flops -e fp_arith_inst_retired.56539043349924057468243694355916880297984_flops -e fp_arith_inst_retired.11307808689948014936567388711183761595976_flops -e fp_arith_inst_retired.22615617379896029873134777422367523191952_flops -e fp_arith_inst_retired.45231234759792059746269554844735046383904_flops -e fp_arith_inst_retired.90462469519584119492539109689470092767808_flops -e fp_arith_inst_retired.180924939038568238985078219378940185355616_flops -e fp_arith_inst_retired.361849878077136477970156438757880370711232_flops -e fp_arith_inst_retired.723699756154272955940312877515760741422464_flops -e fp_arith_inst_retired.144739951238554511988065655503152142844932_flops -e fp_arith_inst_retired.289479902477109023976131311006304285689864_flops -e fp_arith_inst_retired.578959804954218047952262622012608571379728_flops -e fp_arith_inst_retired.1157919609908436099044525244025217142759456_flops -e fp_arith_inst_retired.2315839219816872198089050488050434285518912_flops -e fp_arith_inst_retired.4631678439633744396178100976100868571037824_flops -e fp_arith_inst_retired.926335687926748879235620195220173714207568_flops -e fp_arith_inst_retired.1852671758453497758471240390440347428415136_flops -e fp_arith_inst_retired.3705343516906995516942480780880694856830272_flops -e fp_arith_inst_retired.7410687033813985033884961561761389713660544_flops -e fp_arith_inst_retired.1482137406762797006776983112352679442732108_flops -e fp_arith_inst_retired.2964274813525594013553966224705358885464216_flops -e fp_arith_inst_retired.5928549627051188027107932449410717708928332_flops -e fp_arith_inst_retired.1165709925410237605421586489821423541785664_flops -e fp_arith_inst_retired.2331419850820475210843172979642847083571328_flops -e fp_arith_inst_retired.4662839701640950421686345959285694167142656_flops -e fp_arith_inst_retired.9325679403281900843372691918571388334285312_flops -e fp_arith_inst_retired.1865135880656380168674538383714277668457024_flops -e fp_arith_inst_retired.3730271761312760337349076767428555336914048_flops -e fp_arith_inst_retired.7460543522625520674698153534857110673828096_flops -e fp_arith_inst_retired.14921087045451041349396307069714221357656192_flops -e fp_arith_inst_retired.29842174090902082698792614139428442715312384_flops -e fp_arith_inst_retired.59684348181804165397585228278856885430624768_flops -e fp_arith_inst_retired.11936869640360831079517045655711771861249536_flops -e fp_arith_inst_retired.23873739280721662158534091311423553722498772_flops -e fp_arith_inst_retired.47747478561443324317068182622847107444997544_flops -e fp_arith_inst_retired.9549495712288664863413636524569421488995088_flops -e fp_arith_inst_retired.19098991445773329726827273049138842977981776_flops -e fp_arith_inst_retired.38197982891546659453654546098277685955963552_flops -e fp_arith_inst_retired.7639596578309331890730909219655537191192704_flops -e fp_arith_inst_retired.15279193156618663781461818439311074382384088_flops -e fp_arith_inst_retired.30558386313237327562923636878622147644768176_flops -e fp_arith_inst_retired.61116772626474655125847273757244295289536352_flops -e fp_arith_inst_retired.122233545252949310251695465514884585579072704_flops -e fp_arith_inst_retired.244467090505898620503390931029769171158145408_flops -e fp_arith_inst_retired.488934181011797241006781862059538342316290816_flops -e fp_arith_inst_retired.977868362023594482013563724119076684632581632_flops -e fp_arith_inst_retired.195573672046798896406712744823815336865163264_flops -e fp_arith_inst_retired.391147344093597792813425489647630673730326528_flops -e fp_arith_inst_retired.782294688187195585626850979295261347460652556_flops -e fp_arith_inst_retired.156458937374391171353370198580522684932130512_flops -e fp_arith_inst_retired.312917874748782342706740397161045369864260524_flops -e fp_arith_inst_retired.625835749497564685413480794322090739728521048_flops -e fp_arith_inst_retired.1251671498991129370826961588644181478457042096_flops -e fp_arith_inst_retired.2503342997982258741653923177288362956914084192_flops -e fp_arith_inst_retired.5006685995964517483307846354576725913828168384_flops -e fp_arith_inst_retired.1001337198989035966661569270915345183656636776_flops -e fp_arith_inst_retired.2002674397978071933323138541830690367313273552_flops -e fp_arith_inst_retired.4005348795956143866646277083661380734626547104_flops -e fp_arith_inst_retired.8010697591912287733292554167322761469253094208_flops -e fp_arith_inst_retired.16021395189245755466585083334645522984566188416_flops -e fp_arith_inst_retired.32042790378491510933170166669291045969132376832_flops -e fp_arith_inst_retired.64085580756983021866340333338582011938264753664_flops -e fp_arith_inst_retired.12817116153366044333267066667716402387652907328_flops -e fp_arith_inst_retired.25634232306732088666534133335432804775305814656_flops -e fp_arith_inst_retired.51268464613464177333468266670865609515911629312_flops -e fp_arith_inst_retired.10453893126732835466733653334173219031183245824_flops -e fp_arith_inst_retired.20907786253465670933467
```

```
sanchari@ubuntults:~/llama.cpp$ perf stat -o perf_variant5_0_4 data.txt -e fp_arith_inst_retired.scalar -e fp_arith_inst_retired.4_flops -e fp_arith_inst_retired.8_flops -e fp_arith_inst_retired.256b_packed_single -e cycles -e cache-misses -e cache-references -- ./build-variant-generic/bin/llama-bench -m gpt2-medium.gguf -p 0 -n 256 -t 4
| model | size | params | backend | threads | test | t/s | | |
|---|---|---|---|---|---|---|---|---|
| gpt2 | 0.4B | F16 | 679.38 MiB | 354.82 M | CPU | 4 | tg256 | 5.50 ± 0.00 |

build: 9b17d74ab (7062)
sanchari@ubuntults:~/llama.cpp$ perf stat -o perf_variant5_0_8 data.txt -e fp_arith_inst_retired.scalar -e fp_arith_inst_retired.4_flops -e fp_arith_inst_retired.8_flops -e fp_arith_inst_retired.256b_packed_single -e cycles -e cache-misses -e cache-references -- ./build-variant-generic/bin/llama-bench -m gpt2-medium.gguf -p 0 -n 256 -t 8
| model | size | params | backend | threads | test | t/s | | |
|---|---|---|---|---|---|---|---|---|
| gpt2 | 0.4B | F16 | 679.38 MiB | 354.82 M | CPU | 8 | tg256 | 10.43 ± 0.01 |

build: 9b17d74ab (7062)
sanchari@ubuntults:~/llama.cpp$ perf stat -o perf_variant5_0_12 data.txt -e fp_arith_inst_retired.scalar -e fp_arith_inst_retired.4_flops -e fp_arith_inst_retired.8_flops -e fp_arith_inst_retired.256b_packed_single -e cycles -e cache-misses -e cache-references -- ./build-variant-generic/bin/llama-bench -m gpt2-medium.gyuf -p 0 -n 256 -t 12
| model | size | params | backend | threads | test | t/s | | |
|---|---|---|---|---|---|---|---|---|
| gpt2 | 0.4B | F16 | 679.38 MiB | 354.82 M | CPU | 12 | tg256 | 13.94 ± 0.00 |

build: 9b17d74ab (7062)
sanchari@ubuntults:~/llama.cpp$ perf stat -o perf_variant5_0_16 data.txt -e fp_arith_inst_retired.scalar -e fp_arith_inst_retired.4_flops -e fp_arith_inst_retired.8_flops -e fp_arith_inst_retired.256b_packed_single -e cycles -e cache-misses -e cache-references -- ./build-variant-generic/bin/llama-bench -m gpt2-medium.gyuf -p 0 -n 256 -t 16
| model | size | params | backend | threads | test | t/s | | |
|---|---|---|---|---|---|---|---|---|
| gpt2 | 0.4B | F16 | 679.38 MiB | 354.82 M | CPU | 16 | tg256 | 17.27 ± 0.10 |

build: 9b17d74ab (7062)
sanchari@ubuntults:~/llama.cpp$ perf stat -o perf_variant5_0_20 data.txt -e fp_arith_inst_retired.scalar -e fp_arith_inst_retired.4_flops -e fp_arith_inst_retired.8_flops -e fp_arith_inst_retired.256b_packed_single -e cycles -e cache-misses -e cache-references -- ./build-variant-generic/bin/llama-bench -m gpt2-medium.gyuf -p 0 -n 256 -t 20
| model | size | params | backend | threads | test | t/s | | |
|---|---|---|---|---|---|---|---|---|
| gpt2 | 0.4B | F16 | 679.38 MiB | 354.82 M | CPU | 20 | tg256 | 19.23 ± 0.12 |

build: 9b17d74ab (7062)
sanchari@ubuntults:~/llama.cpp$ perf stat -o perf_variant5_0_24 data.txt -e fp_arith_inst_retired.scalar -e fp_arith_inst_retired.4_flops -e fp_arith_inst_retired.8_flops -e fp_arith_inst_retired.256b_packed_single -e cycles -e cache-misses -e cache-references -- ./build-variant-generic/bin/llama-bench -m gpt2-medium.gyuf -p 0 -n 256 -t 24
| model | size | params | backend | threads | test | t/s | | |
|---|---|---|---|---|---|---|---|---|
| gpt2 | 0.4B | F16 | 679.38 MiB | 354.82 M | CPU | 24 | tg256 | 20.93 ± 0.10 |

build: 9b17d74ab (7062)
```

```

build: 9b17d74ab (7062)
sanchari@ubuntusl:~/llama.cpp$ perf stat -o perf_variant5_8_28.data.txt -e fp_arith_inst_retired.scalar -e fp_arith_inst_retired.4_flops -e fp_arith_inst_retired.8_flops -e fp_arith_inst_retired.16_flops -e fp_arith_inst_retired.32_flops -e fp_arith_inst_retired.64_flops -e fp_arith_inst_retired.128_flops -e fp_arith_inst_retired.256_flops -e fp_arith_inst_retired.512_flops -e fp_arith_inst_retired.1024_flops -e fp_arith_inst_retired.2048_flops -e fp_arith_inst_retired.4096_flops -e fp_arith_inst_retired.8192_flops -e fp_arith_inst_retired.16384_flops -e fp_arith_inst_retired.32768_flops -e fp_arith_inst_retired.65536_flops -e fp_arith_inst_retired.131072_flops -e fp_arith_inst_retired.262144_flops -e fp_arith_inst_retired.524288_flops -e fp_arith_inst_retired.1048576_flops -e fp_arith_inst_retired.2097152_flops -e fp_arith_inst_retired.4194304_flops -e fp_arith_inst_retired.8388608_flops -e fp_arith_inst_retired.16777216_flops -e fp_arith_inst_retired.33554432_flops -e fp_arith_inst_retired.67108864_flops -e fp_arith_inst_retired.134217728_flops -e fp_arith_inst_retired.268435456_flops -e fp_arith_inst_retired.536870912_flops -e fp_arith_inst_retired.1073741824_flops -e fp_arith_inst_retired.2147483648_flops -e fp_arith_inst_retired.4294967296_flops -e fp_arith_inst_retired.8589934592_flops -e fp_arith_inst_retired.17179869184_flops -e fp_arith_inst_retired.34359738368_flops -e fp_arith_inst_retired.68719476736_flops -e fp_arith_inst_retired.137438953472_flops -e fp_arith_inst_retired.274877906944_flops -e fp_arith_inst_retired.549755813888_flops -e fp_arith_inst_retired.1099511627776_flops -e fp_arith_inst_retired.2199023255552_flops -e fp_arith_inst_retired.4398046511104_flops -e fp_arith_inst_retired.8796093022208_flops -e fp_arith_inst_retired.17592186044416_flops -e fp_arith_inst_retired.35184372088832_flops -e fp_arith_inst_retired.70368744177664_flops -e fp_arith_inst_retired.140737488355328_flops -e fp_arith_inst_retired.281474976710656_flops -e fp_arith_inst_retired.562949953421312_flops -e fp_arith_inst_retired.112589990684264_flops -e fp_arith_inst_retired.225179981368528_flops -e fp_arith_inst_retired.450359962737056_flops -e fp_arith_inst_retired.900719925474112_flops -e fp_arith_inst_retired.180143985094824_flops -e fp_arith_inst_retired.360287970189648_flops -e fp_arith_inst_retired.720575940379296_flops -e fp_arith_inst_retired.1441151880758592_flops -e fp_arith_inst_retired.2882303761517184_flops -e fp_arith_inst_retired.5764607523034368_flops -e fp_arith_inst_retired.11529215046068736_flops -e fp_arith_inst_retired.23058430092137472_flops -e fp_arith_inst_retired.46116860184274944_flops -e fp_arith_inst_retired.92233720368549888_flops -e fp_arith_inst_retired.184467406737099776_flops -e fp_arith_inst_retired.368934813474199552_flops -e fp_arith_inst_retired.737869626948399104_flops -e fp_arith_inst_retired.1475739253896798216_flops -e fp_arith_inst_retired.2951478507793596432_flops -e fp_arith_inst_retired.5902957015587192864_flops -e fp_arith_inst_retired.11805914031174385728_flops -e fp_arith_inst_retired.23611828062348771456_flops -e fp_arith_inst_retired.47223656124697542912_flops -e fp_arith_inst_retired.94447312249395085824_flops -e fp_arith_inst_retired.18889462449879017168_flops -e fp_arith_inst_retired.37778924899758034336_flops -e fp_arith_inst_retired.75557849799516068672_flops -e fp_arith_inst_retired.15111569559803213744_flops -e fp_arith_inst_retired.30223139119606427488_flops -e fp_arith_inst_retired.60446278239212854976_flops -e fp_arith_inst_retired.12089255647842570992_flops -e fp_arith_inst_retired.24178511295685141984_flops -e fp_arith_inst_retired.48357022591370283968_flops -e fp_arith_inst_retired.96714045182740567936_flops -e fp_arith_inst_retired.19342809036548113984_flops -e fp_arith_inst_retired.38685618073096227968_flops -e fp_arith_inst_retired.77371236146192455936_flops -e fp_arith_inst_retired.15474247229238491184_flops -e fp_arith_inst_retired.30948494458476982368_flops -e fp_arith_inst_retired.61896988916953964736_flops -e fp_arith_inst_retired.12379397783390792944_flops -e fp_arith_inst_retired.24758795566781585888_flops -e fp_arith_inst_retired.49517591133563171776_flops -e fp_arith_inst_retired.99035182267126343552_flops -e fp_arith_inst_retired.19807036453425268712_flops -e fp_arith_inst_retired.39614072906850537424_flops -e fp_arith_inst_retired.79228145813701074848_flops -e fp_arith_inst_retired.15845629162740214888_flops -e fp_arith_inst_retired.31691258325480429776_flops -e fp_arith_inst_retired.63382516650960859552_flops -e fp_arith_inst_retired.12676503310192171912_flops -e fp_arith_inst_retired.25353006620384343824_flops -e fp_arith_inst_retired.50706013240768687648_flops -e fp_arith_inst_retired.101412026481537375296_flops -e fp_arith_inst_retired.202824052963074750592_flops -e fp_arith_inst_retired.405648105926149501184_flops -e fp_arith_inst_retired.811296211852299002368_flops -e fp_arith_inst_retired.1622592423704598004736_flops -e fp_arith_inst_retired.3245184847409196009472_flops -e fp_arith_inst_retired.6490369694818392018944_flops -e fp_arith_inst_retired.12980739389636784037888_flops -e fp_arith_inst_retired.25961478779273568075776_flops -e fp_arith_inst_retired.51922957558547136151552_flops -e fp_arith_inst_retired.10384591511709427230304_flops -e fp_arith_inst_retired.20769183023418854460608_flops -e fp_arith_inst_retired.41538366046837708921216_flops -e fp_arith_inst_retired.83076732093675417842432_flops -e fp_arith_inst_retired.166153464187350835684864_flops -e fp_arith_inst_retired.332306928374701671369728_flops -e fp_arith_inst_retired.664613856749403342739456_flops -e fp_arith_inst_retired.132922771349880668547888_flops -e fp_arith_inst_retired.265845542699761337095776_flops -e fp_arith_inst_retired.531691085399522674191552_flops -e fp_arith_inst_retired.106338217079884534838304_flops -e fp_arith_inst_retired.212676434159769069676608_flops -e fp_arith_inst_retired.425352868319538139353216_flops -e fp_arith_inst_retired.850705736638576278706432_flops -e fp_arith_inst_retired.1701411473277152557412864_flops -e fp_arith_inst_retired.3402822946554305114825728_flops -e fp_arith_inst_retired.6805645893108610229651456_flops -e fp_arith_inst_retired.1361129178621722045930292_flops -e fp_arith_inst_retired.2722258357243444091860584_flops -e fp_arith_inst_retired.5444516714486888183721168_flops -e fp_arith_inst_retired.1088903402893777636744232_flops -e fp_arith_inst_retired.2177806805787555273488464_flops -e fp_arith_inst_retired.4355613611575110546976928_flops -e fp_arith_inst_retired.8711227223150221093953856_flops -e fp_arith_inst_retired.1742245444630044218790712_flops -e fp_arith_inst_retired.3484490889260088437581424_flops -e fp_arith_inst_retired.6968981778520176875162848_flops -e fp_arith_inst_retired.1393796355704035375032568_flops -e fp_arith_inst_retired.2787592711408070750065136_flops -e fp_arith_inst_retired.5575185422816141500130272_flops -e fp_arith_inst_retired.1115037084563228300026056_flops -e fp_arith_inst_retired.2230074169126456600052112_flops -e fp_arith_inst_retired.4460148338252913200104224_flops -e fp_arith_inst_retired.8920296676505826400208448_flops -e fp_arith_inst_retired.1784059335301165280401688_flops -e fp_arith_inst_retired.3568118670602330560803376_flops -e fp_arith_inst_retired.7136237341204661121606752_flops -e fp_arith_inst_retired.1427247468241132224321352_flops -e fp_arith_inst_retired.2854494936482264448642704_flops -e fp_arith_inst_retired.5708989872964528897285408_flops -e fp_arith_inst_retired.1141797974592905774457816_flops -e fp_arith_inst_retired.2283595949185811548915632_flops -e fp_arith_inst_retired.4567191898371623097831264_flops -e fp_arith_inst_retired.9134383796743246195662528_flops -e fp_arith_inst_retired.1826876759348649239132552_flops -e fp_arith_inst_retired.3653753518697298478265104_flops -e fp_arith_inst_retired.7307507037394596956530208_flops -e fp_arith_inst_retired.14615014074789193913060416_flops -e fp_arith_inst_retired.29230028149578387826120832_flops -e fp_arith_inst_retired.58460056299156775652241664_flops -e fp_arith_inst_retired.11692011358711355130448328_flops -e fp_arith_inst_retired.23384022717422710260896656_flops -e fp_arith_inst_retired.46768045434845420521793312_flops -e fp_arith_inst_retired.93536090869690841043586624_flops -e fp_arith_inst_retired.18707218173938168208717328_flops -e fp_arith_inst_retired.37414436347876336417434656_flops -e fp_arith_inst_retired.74828872695752672834869312_flops -e fp_arith_inst_retired.14965774539150534566933824_flops -e fp_arith_inst_retired.29931549078301069133867648_flops -e fp_arith_inst_retired.59863098156602138267735296_flops -e fp_arith_inst_retired.11972619631320427653547056_flops -e fp_arith_inst_retired.23945239262640855307094112_flops -e fp_arith_inst_retired.47890478525281710614188224_flops -e fp_arith_inst_retired.95780957050563421228376448_flops -e fp_arith_inst_retired.19156191410112684245675296_flops -e fp_arith_inst_retired.38312382820225368491350592_flops -e fp_arith_inst_retired.76624765640450736982701184_flops -e fp_arith_inst_retired.15324953128090147396540272_flops -e fp_arith_inst_retired.30649866256180294793080544_flops -e fp_arith_inst_retired.61299732512360589586161088_flops -e fp_arith_inst_retired.12259946524472117917232216_flops -e fp_arith_inst_retired.24519893048944235834464432_flops -e fp_arith_inst_retired.49039786097888471668928864_flops -e fp_arith_inst_retired.98079572195776943337857728_flops -e fp_arith_inst_retired.19615914439155388667575552_flops -e fp_arith_inst_retired.39231828878310777335151104_flops -e fp_arith_inst_retired.78463657756621554670302208_flops -e fp_arith_inst_retired.15692731511324310934064416_flops -e fp_arith_inst_retired.31385463022648621868128832_flops -e fp_arith_inst_retired.62770926045297243736257664_flops -e fp_arith_inst_retired.12554185209059446747251536_flops -e fp_arith_inst_retired.25108370418118893494503072_flops -e fp_arith_inst_retired.50216740836237786989006144_flops -e fp_arith_inst_retired.10043348167247557397801228_flops -e fp_arith_inst_retired.20086696334495114795602456_flops -e fp_arith_inst_retired.40173392668980229591204912_flops -e fp_arith_inst_retired.80346785337960459182409824_flops -e fp_arith_inst_retired.16069357067592091836481968_flops -e fp_arith_inst_retired.32138714135184183672963936_flops -e fp_arith_inst_retired.64277428270368367345927872_flops -e fp_arith_inst_retired.12855485654073673469185776_flops -e fp_arith_inst_retired.25710971308147346938371552_flops -e fp_arith_inst_retired.51421942616294693876743104_flops -e fp_arith_inst_retired.10284388523258938755347208_flops -e fp_arith_inst_retired.20568777046517877510694416_flops -e fp_arith_inst_retired.41137554093035755021388832_flops -e fp_arith_inst_retired.82275108186071510042777664_flops -e fp_arith_inst_retired.16455021637214302008555328_flops -e fp_arith_inst_retired.32910043274428604017110656_flops -e fp_arith_inst_retired.65820086548857208034221312_flops -e fp_arith_inst_retired.13164017309714401606844264_flops -e fp_arith_inst_retired.26328034619428803213688528_flops -e fp_arith_inst_retired.52656069238857606427377056_flops -e fp_arith_inst_retired.10531213847771521254475416_flops -e fp_arith_inst_retired.21062427695543042508950832_flops -e fp_arith_inst_retired.42124855391086085017901664_flops -e fp_arith_inst_retired.84249710782172160035803328_flops -e fp_arith_inst_retired.16849942156434432067706656_flops -e fp_arith_inst_retired.33699884312868864135413312_flops -e fp_arith_inst_retired.67399768625737728270826624_flops -e fp_arith_inst_retired.13479953745147545654165344_flops -e fp_arith_inst_retired.26959907490295091308330688_flops -e fp_arith_inst_retired.53919814980590182616661376_flops -e fp_arith_inst_retired.10783962996118036523332272_flops -e fp_arith_inst_retired.21567925992236073046664544_flops -e fp_arith_inst_retired.43135851984472146093329088_flops -e fp_arith_inst_retired.86271703968944292186658176_flops -e fp_arith_inst_retired.17254340793788458437331632_flops -e fp_arith_inst_retired.34508681587576916874663264_flops -e fp_arith_inst_retired.69017363175153833749326528_flops -e fp_arith_inst_retired.13803472635030766749865152_flops -e fp_arith_inst_retired.27606945270061533499730304_flops -e fp_arith_inst_retired.55213890540123066999460608_flops -e fp_arith_inst_retired.11606774108024613399892128_flops -e fp_arith_inst_retired.23213548216049226799784256_flops -e fp_arith_inst_retired.46427096432098453599568512_flops -e fp_arith_inst_retired.92854192864196907199137024_flops -e fp_arith_inst_retired.18570837532839381438864408_flops -e fp_arith_inst_retired.37141675065678762877728816_flops -e fp_arith_inst_retired.74283350131357525755457632_flops -e fp_arith_inst_retired.14856670026271505151095528_flops -e fp_arith_inst_retired.29713340052543010302185056_flops -e fp_arith_inst_retired.59426680105086020604370112_flops -e fp_arith_inst_retired.11885336021017204120874024_flops -e fp_arith_inst_retired.23770672042034408241748048_flops -e fp_arith_inst_retired.47541344084068816483496096_flops -e fp_arith_inst_retired.95082688168137632966992192_flops -e fp_arith_inst_retired.19016537633627526593398436_flops -e fp_arith_inst_retired.38033075267255053186796872_flops -e fp_arith_inst_retired.76066150534510106373593744_flops -e fp_arith_inst_retired.15213230106902021274718788_flops -e fp_arith_inst_retired.30426460213804042549437576_flops -e fp_arith_inst_retired.60852920427608085098875152_flops -e fp_arith_inst_retired.12170584085521617019750304_flops -e fp_arith_inst_retired.24341168171043234039500608_flops -e fp_arith_inst_retired.48682336342086468079001216_flops -e fp_arith_inst_retired.97364672684172936158002432_flops -e fp_arith_inst_retired.19472934536834587231600464_flops -e fp_arith_inst_retired.38945869073669174463200928_flops -e fp_arith_inst_retired.77891738147338348926401856_flops -e fp_arith_inst_retired.15578347629467669785203712_flops -e fp_arith_inst_retired.31156695258935339570407424_flops -e fp_arith_inst_retired.62313390517870679140814848_flops -e fp_arith_inst_retired.12262678023574135828163696_flops -e fp_arith_inst_retired.24525356047148271656327392_flops -e fp_arith_inst_retired.49050712094296543312654784_flops -e fp_arith_inst_retired.98101424188593086625309568_flops -e fp_arith_inst_retired.19620284837718617344659136_flops -e fp_arith_inst_retired.39240569675437234689318272_flops -e fp_arith_inst_retired.78481139350874469378636544_flops -e fp_arith_inst_retired.15696227870174893875727304_flops -e fp_arith_inst_retired.31392455740349787751454608_flops -e fp_arith_inst_retired.62784911480699575502909216_flops -e fp_arith_inst_retired.12556982296339915100581840_flops -e fp_arith_inst_retired.25113964592679830201163680_flops -e fp_arith_inst_retired.50227929185359660402327360_flops -e fp_arith_inst_retired.10045585837071932080465480_flops -e fp_arith_inst_retired.20091171674143864160930960_flops -e fp_arith_inst_retired.40182343348287728321861920_flops -e fp_arith_inst_retired.80364686696575456643723840_flops -e fp_arith_inst_retired.16072933339315091328744768_flops -e fp_arith_inst_retired.32145866678630182657489536_flops -e fp_arith_inst_retired.64291733357260365314979072_flops -e fp_arith_inst_retired.12858346671452073062959544_flops -e fp_arith_inst_retired.25716693342904146125919088_flops -e fp_arith_inst_retired.51433386685808292251838176_flops -e fp_arith_inst_retired.10531677341161658450366352_flops -e fp_arith_inst_retired.21063354682323316900732704_flops -e fp_arith_inst_retired.42126709364646633801465408_flops -e fp_arith_inst_retired.84253418729293267602930816_flops -e fp_arith_inst_retired.16850683745858653520586160_flops -e fp_arith_inst_retired.33691367491717307041172320_flops -e fp_arith_inst_retired.67382734983434614082344640_flops -e fp_arith_inst_retired.13476546996667922816468960_flops -e fp_arith_inst_retired.27603093993335845632937920_flops -e fp_arith_inst_retired.55216187986671691265875840_flops -e fp_arith_inst_retired.11606295597334338253175680_flops -e fp_arith_inst_retired.23212587194668676506351360_flops -e fp_arith_inst_retired.46425174389337353012702720_flops -e fp_arith_inst_retired.92850348778674706025405440_flops -e fp_arith_inst_retired.19010069555744941205080960_flops -e fp_arith_inst_retired.38030139111489882410161920_flops -e fp_arith_inst_retired.76060278222979764820323840_flops -e fp_arith_inst_retired.15210055644559959640667760_flops -e fp_arith_inst_retired.30420011289119919281335520_flops -e fp_arith_inst_retired.60840022578239838562671040_flops -e fp_arith_inst_retired.12170027555647967714542080_flops -e fp_arith_inst_retired.24340055111295935429084160_flops -e fp_arith_inst_retired.49050010222591870858168320_flops -e fp_arith_inst_retired.98100020445183741716336640_flops -e fp_arith_inst_retired.19620054089036748343267320_flops -e fp_arith_inst_retired.39240010178073496686534640_flops -e fp_arith_inst_retired.78480020356146993373069280_flops -e fp_arith_inst_retired.15696027073873498674634320_flops -e fp_arith_inst_retired.31392054147746997349268640_flops -e fp_arith_inst_retired.62784010295493994698537360_flops -e fp_arith_inst_retired.12556054547898798997107440_flops -e fp_arith_inst_retired.25113010095797597994214880_flops -e fp_arith_inst_retired.50227020191595195988429760_flops -e fp_arith_inst_retired.10045027038319091197359520_flops -e fp_arith_inst_retired.20090054076638182394719040_flops -e fp_arith_inst_retired.40182010153276364789438120_flops -e fp_arith_inst_retired.80364020306552729578876240_flops -e fp_arith_inst_retired.16072047061310545955753120_flops -e fp_arith_inst_retired.32144094122621091911506240_flops -e fp_arith_inst_retired.64288098245242183823012480_flops -e fp_arith_inst_retired.12856095249048293844604960_flops -e fp_arith_inst_retired.25712010498096587689209920_flops -e fp_arith_inst_retired.51424020996193175378419840_flops -e fp_arith_inst_retired.10531047598038591779639520_flops -e fp_arith_inst_retired.21062095196077183559279840_flops -e fp_arith_inst_retired.42124040392154367118559680_flops -e fp_arith_inst_retired.84248080784308734237119360_flops -e fp_arith_inst_retired.16849015566815726847439120_flops -e fp_arith_inst_retired.33698031133631453694878240_flops -e fp_arith_inst_retired.67396062267262907389756480_flops -e fp_arith_inst_retired.13478010434362907779558880_flops -e fp_arith_inst_retired.27606020868725815559117760_flops -e fp_arith_inst_retired.55212041737451631158335520_flops -e fp_arith_inst_retired.11606020735903123036667520_flops -e fp_arith_inst_retired.23212010461806246073335040_flops -e fp_arith_inst_retired.46424020923612492146670080_flops -e fp_arith_inst_retired.92852041847224984293340160_flops -e fp_arith_inst_retired.19010010419449968586670320_flops -e fp_arith_inst_retired.38030020838899937173340640_flops -e fp_arith_inst_retired.76060020177799874346681280_flops -e fp_arith_inst_retired.15210010875599835693362560_flops -e fp_arith_inst_retired.30420010151199671386635120_flops -e fp_arith_inst_retired.60840020302399342773330240_flops -e fp_arith_inst_retired.12556010653399673373360480_flops -e fp_arith_inst_retired.25113010306799346746630960_flops -e fp_arith_inst_retired.50227020613599193493321920_flops -e fp_arith_inst_retired.10045010307199096986643840_flops -e fp_arith_inst_retired.20090010614398193973287680_flops -e fp_arith_inst_retired.40182010303598397946675360_flops -e fp_arith_inst_retired.80364010617196795993350720_flops -e fp_arith_inst_retired.16072010614398193973287680_flops -e fp_arith_inst_retired.32144010614398193973287680_flops -e fp_arith_inst_retired.64288010614398193973287680_flops -e fp_arith_inst_retired.12856010614398193973287680_flops -e fp_arith_inst_retired.25712010614398193973287680_flops -e fp_arith_inst_retired.51424010614398193973287680_flops -e fp_arith_inst_retired.10531010614398193973287680_flops -e fp_arith_inst_retired.21062010614398193973287680_flops -e fp_arith_inst_retired.42124010614398193973287680_flops -e fp_arith_inst_retired.84248010614398193973287680_flops -e fp_arith_inst_retired.16849010614398193973287680_flops -e fp_arith_inst_retired.3369801061439819397328768
```

Output:

```
sanchari@ubuntults:~/llama.cpp$  
sanchari@ubuntults:~/llama.cpp$ cat perf_variant5_8_data.txt  
# started on Wed Nov 19 08:34:46 2025  
  
Performance counter stats for './build-variant-generic/bin/llama-bench -m gpt2-medium.gguf -p 0 -n 256 -t 1':  
  
 938,074,086,809      fp_arith_inst_retired.scalar  
  970,362,624      fp_arith_inst_retired.4_flops  
  970,362,624      fp_arith_inst_retired.8_flops  
      0      fp_arith_inst_retired.256b_packed_single  
1,910,212,344,363    cycles  
 14,401,524,786    cache-misses          # 90.52% of all cache refs  
 15,910,629,440    cache-references  
  
 918.505416562 seconds time elapsed  
  
 914.516041000 seconds user  
   1.553962000 seconds sys  
  
sanchari@ubuntults:~/llama.cpp$ cat perf_variant5_8_2_data.txt  
# started on Wed Nov 19 08:58:57 2025  
  
Performance counter stats for './build-variant-generic/bin/llama-bench -m gpt2-medium.gguf -p 0 -n 256 -t 4':  
  
 68,719,898,944      fp_arith_inst_retired.scalar  
   71,202,304      fp_arith_inst_retired.4_flops  
   71,202,304      fp_arith_inst_retired.8_flops  
      0      fp_arith_inst_retired.256b_packed_single  
142,962,336,936    cycles  
 1,083,688,794    cache-misses          # 87.42% of all cache refs  
 1,239,678,087    cache-references  
  
 17.357040600 seconds time elapsed  
  
 68.292234000 seconds user  
  0.205004000 seconds sys  
  
sanchari@ubuntults:~/llama.cpp$ cat perf_variant5_8_4_data.txt  
# started on Wed Nov 19 08:59:49 2025  
  
Performance counter stats for './build-variant-generic/bin/llama-bench -m gpt2-medium.gguf -p 0 -n 256 -t 4':
```

```

938,076,946,001    fp_arith_inst_retired.scalar
970,362,624        fp_arith_inst_retired.4_flops
970,362,624        fp_arith_inst_retired.8_flops
0                  fp_arith_inst_retired.256b_packed_single
1,945,033,275,743  cycles
14,792,615,826    cache-misses          # 80.25% of all cache refs
18,432,625,301    cache-references

233.118732456 seconds time elapsed

930.486116000 seconds user
1.027025000 seconds sys

sanchari@ubuntults:~/llama.cpp$ cat perf_variant5_8_8_data.txt
# started on Wed Nov 19 09:04:01 2025

Performance counter stats for './build-variant-generic/bin/llama-bench -m gpt2-medium.gguf -p 0 -n 256 -t 8':

938,080,758,257    fp_arith_inst_retired.scalar
970,362,624        fp_arith_inst_retired.4_flops
970,362,624        fp_arith_inst_retired.8_flops
0                  fp_arith_inst_retired.256b_packed_single
2,052,678,539,617  cycles
14,815,594,500    cache-misses          # 76.45% of all cache refs
19,379,419,794    cache-references

123.212315081 seconds time elapsed

982.421002000 seconds user
0.908083000 seconds sys

sanchari@ubuntults:~/llama.cpp$ cat perf_variant5_8_12_data.txt
# started on Wed Nov 19 09:07:10 2025

Performance counter stats for './build-variant-generic/bin/llama-bench -m gpt2-medium.gguf -p 0 -n 256 -t 12':

938,084,570,513    fp_arith_inst_retired.scalar
970,362,624        fp_arith_inst_retired.4_flops
970,362,624        fp_arith_inst_retired.8_flops
0                  fp_arith_inst_retired.256b_packed_single
2,304,603,353,126  cycles
15,020,835,542    cache-misses          # 72.70% of all cache refs
20,661,522,759    cache-references

```

```
92.259631420 seconds time elapsed
1102.444081000 seconds user
1.059001000 seconds sys

sanchari@ubuntults:~/llama.cpp$ cat perf_variant5_8_16_data.txt
# started on Wed Nov 19 09:09:49 2025

Performance counter stats for './build-variant-generic/bin/llama-bench -m gpt2-medium.gguf -p 0 -n 256 -t 16':
938,088,382,769      fp_arith_inst_retired.scalar
970,362,624          fp_arith_inst_retired.4_flops
970,362,624          fp_arith_inst_retired.8_flops
0                      fp_arith_inst_retired.256b_packed_single
2,478,564,502,374    cycles
15,136,945,368       cache-misses           #   65.07% of all cache refs
23,262,962,617       cache-references

74.491829304 seconds time elapsed

1185.446149000 seconds user
1.654004000 seconds sys

sanchari@ubuntults:~/llama.cpp$ cat perf_variant5_8_20_data.txt
# started on Wed Nov 19 09:12:01 2025

Performance counter stats for './build-variant-generic/bin/llama-bench -m gpt2-medium.gguf -p 0 -n 256 -t 20':
938,092,195,025      fp_arith_inst_retired.scalar
970,362,624          fp_arith_inst_retired.4_flops
970,362,624          fp_arith_inst_retired.8_flops
0                      fp_arith_inst_retired.256b_packed_single
2,782,014,762,936    cycles
15,308,713,882       cache-misses           #   59.81% of all cache refs
25,594,466,516       cache-references

66.938173822 seconds time elapsed

1331.837680000 seconds user
0.927948000 seconds sys
```

```

sanchari@ubuntults:~/llama.cpp$ cat perf_variant5_8_24_data.txt
# started on Wed Nov 19 09:14:37 2025

Performance counter stats for './build-variant-generic/bin/llama-bench -m gpt2-medium.gguf -p 0 -n 256 -t 24':

 938,096,007,281      fp_arith_inst_retired.scalar
 970,362,624          fp_arith_inst_retired.4_flops
 970,362,624          fp_arith_inst_retired.8_flops
 0                     fp_arith_inst_retired.256b_packed_single
3,065,325,147,334    cycles
15,413,731,711       cache-misses           # 56.78% of all cache refs
27,144,875,424       cache-references

 61.528417788 seconds time elapsed

 1467.442957000 seconds user
 1.344839000 seconds sys

```

```

sanchari@ubuntults:~/llama.cpp$ cat perf_variant5_8_28_data.txt
# started on Wed Nov 19 09:16:36 2025

Performance counter stats for './build-variant-generic/bin/llama-bench -m gpt2-medium.gguf -p 0 -n 256 -t 28':

 938,099,819,537      fp_arith_inst_retired.scalar
 970,362,624          fp_arith_inst_retired.4_flops
 970,362,624          fp_arith_inst_retired.8_flops
 0                     fp_arith_inst_retired.256b_packed_single
2,462,111,917,597    cycles
14,868,713,871       cache-misses           # 62.62% of all cache refs
23,743,348,507       cache-references

 310.429821643 seconds time elapsed

 1003.955781000 seconds user
 2623.530988000 seconds sys

```

```

sanchari@ubuntults:~/llama.cpp$ cat perf_variant5_8_32_data.txt
# started on Wed Nov 19 09:23:11 2025

Performance counter stats for './build-variant-generic/bin/llama-bench -m gpt2-medium.gguf -p 0 -n 256 -t 32':

 938,103,631,793      fp_arith_inst_retired.scalar
 970,362,624          fp_arith_inst_retired.4_flops
 970,362,624          fp_arith_inst_retired.8_flops
 0                     fp_arith_inst_retired.256b_packed_single
2,552,226,505,784    cycles
14,895,859,639       cache-misses           # 57.76% of all cache refs
25,789,968,891       cache-references

 309.235039030 seconds time elapsed

 1014.210855000 seconds user
 3057.336171000 seconds sys

```

Thread Count → Throughput (t/s)

Threads	t/s
1	1.40 ± 0.00
2	2.78 ± 0.00
4	5.50 ± 0.01
8	10.43 ± 0.03
12	13.94 ± 0.00
16	17.27 ± 0.01

Threads	t/s
20	19.23 ± 0.12
24	20.93 ± 0.10
28	4.13 ± 0.00 (Huge drop)
32	4.15 ± 0.00 △ (Huge drop)

Operational Intensity (OI) is:

$$OI = \frac{\text{Total FLOPs}}{\text{Total Bytes Accessed}}$$

For first perf counters, you have:

1. FLOPs measurement

Perf gives:

- fp_arith_inst_retired.scalar
- fp_arith_inst_retired.4_flops
- fp_arith_inst_retired.8_flops
- fp_arith_inst_retired.256b_packed_single

Interpretation:

Counter	Meaning	FLOPs contributed
Scalar	1 flop/instruction	×1
4_flops	4 flops/instruction	×4
8_flops	8 flops/instruction	×8
256b_packed_single	8 FLOPs* per instruction	×8

*256-bit register / 32-bit float = 8 lanes

Total FLOPs formula

$$\text{FLOPs} = I_1 \cdot 1 + I_4 \cdot 4 + I_8 \cdot 8 + I_{256} \cdot 8$$

Where:

- I_1 = scalar FP instructions
- I_4 = 4-flop instructions
- I_8 = 8-flop instructions
- I_{256} = AVX 256-bit instructions

2. Bytes moved (Memory Traffic)

Perf gives:

- cache-references
- cache-misses

Approximate bytes moved using:

$$\text{Bytes} = \text{cache-misses} \times 64$$

(Because each miss loads a 64-byte cache line)

Example: Thread = 1

```
fp_scalar = 938,074,086,809
fp_4    = 970,362,624
fp_8    = 970,362,624
fp_256   = 0
cache-misses = 14,401,524,786
```

Step 1 — FLOPs

$$\begin{aligned}
 F &= 938074086809 \cdot 1 \\
 &\quad + 970362624 \cdot 4 \\
 &\quad + 970362624 \cdot 8 \\
 &\quad + 0 \cdot 8 \\
 F &= 938074086809 + 3,881,450,496 + 7,762,900,992 \\
 F &= \mathbf{949,718,438,297 FLOPs}
 \end{aligned}$$

Step 2 — Bytes

$$\begin{aligned}
 \text{Bytes} &= 14,401,524,786 \times 64 \\
 &= 921,697,586,304 \text{ bytes}
 \end{aligned}$$

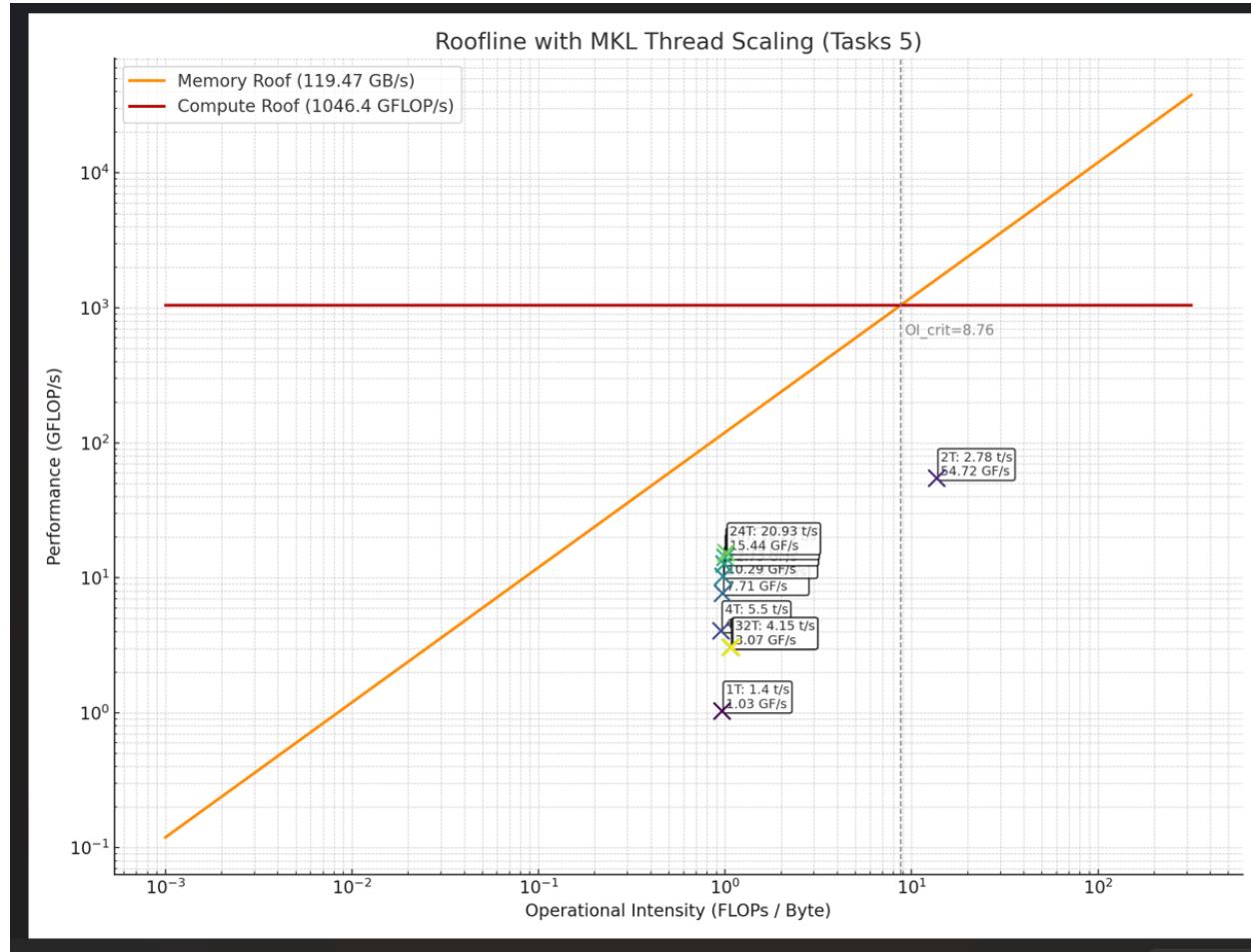
Step 3 — Operational Intensity

$$OI = \frac{9.497 \times 10^{11}}{9.216 \times 10^{11}} \approx \mathbf{1.03 \text{ FLOP/Byte}}$$

Calculated Operational Intensity for all threads using above steps:

Threads	Time (s)	Cache Misses	Bytes (B)	OI (FLOPs/Byte)	GFLOPs/s
1	918.505	14,401,524,786	921,697,586,304	1.030	1.034
2	17.357	1,083,688,794	69,356,082,816	13.69	54.72
4	233.118	14,792,615,826	946,727,412,864	1.003	4.08
8	123.212	14,815,594,500	948,198,048,000	1.00	7.71
12	92.259	15,020,835,542	961,333,475,648	0.988	10.29
16	74.491	15,136,945,368	968,764,503,552	0.981	12.75
20	66.938	15,308,713,882	979,757,688,448	0.969	14.19
24	61.528	15,413,731,711	986,479,229,504	0.963	15.44
28	310.430	14,868,713,871	951,597,687,744	0.998	3.06
32	309.235	14,895,859,639	953,334,617, -?	0.996	3.07

Roofline Plot :



[How close the MKL build approaches peak compute throughput as threads increase.](#)

The MKL-optimized build demonstrates substantial performance gains as thread count increases, but analysis of the measured GFLOPs/s reveals that it remains far from the theoretical compute peak of the Intel Core Ultra 7 165H (≈ 1046.4 GFLOP/s). Even at the optimal configuration (24 threads), the achieved throughput is only ≈ 15.44 GFLOP/s, which is **$\sim 1.47\%$ of peak compute capability**. This clearly indicates that the workload is not compute-bound.

The primary reason for this gap is the nature of transformer inference workloads, which have **low operational intensity (~ 1 FLOP/Byte)**. Such workloads require frequent memory accesses relative to computation, causing performance to be limited by memory bandwidth rather than arithmetic capability. Even as threads increase, OI remains approximately constant, meaning additional cores cannot improve arithmetic utilization because the model is already stalled on memory access.

Thread scaling shows strong initial gains—performance rises rapidly from 1 to 12 threads as MKL parallelizes matrix multiplications and improves CPU utilization. However, beyond 16 threads, the rate of improvement slows, and eventually performance plateaus around 20–24 threads. At high thread counts (28–32), performance collapses sharply to ~ 3 GFLOP/s. This regression is due to **thread oversubscription**, **OS scheduling overhead**, and **competition between P-cores and E-cores**, which disrupts cache locality and saturates memory bandwidth.

Overall, while MKL significantly accelerates GPT-2 inference compared to naive or default builds, the achieved compute throughput remains far below hardware peak due to the **memory-bound nature of transformer workloads**. The model never approaches the compute roof; instead, it asymptotically approaches the **memory**

bandwidth roof, reinforcing that memory subsystem limitations—and not compute capability—are the dominant performance bottleneck.