

# Performance of LLM backends and models in Curnagl

## Introduction

This page shows performance of Llama and mistral models on Curnagl hardware. We have measured the token throughput which should help you to have an idea of what is possible using Curnagl resources. Training time and inference time for different task could be estimated using these results.

---

## Models and backends tested

### Tested Models

#### Llama3

- Official access to Meta Llama3 models: [Meta Llama3 models on Hugging Face](#)
- [Meta-Llama-3.1-8B-Instruct](#)
- [Meta-Llama-3.1-70B-Instruct](#)

#### Mistral

- Official access to Mistral models: [Mistral models on MistralAI website](#)
- Access to Mistral models on Hugging Face: [Mistral models on Hugging Face](#)
- [mistral-7B-Instruct-v0.3](#)
- [Mixtral-8x7B-v0.1-Instruct](#)

---

# Tested Backends

- [vLLM](#)

vLLM backend provides efficient memory usage and fast token sampling. This backend is ideal for testing Llama3 and Mistral models in environments that require high-speed responses and low latency.

- [llama.cpp](#)

llama.cpp was primarily used for llama but it can be applied to other LLM models. This optimized backend provides efficient inference on GPUs.

- [Transformers](#)

If not the most widely used LLM black box, it is one of them. Easy to use, the Hugging Face Transformers library supports a wide range of models and backends. One of its main advantages is its quick set up, which enables quick experimentation across architectures.

- [mistral-inference](#)

This is the official inference backend for Mistral. It is (supposed to be) optimized for Mistral's architecture, thus increasing the model performance. However, our benchmarks results do not demonstrate any specificities to Mistral model as llama.cpp seems to perform better.

---

# Hardware description

Three different types of GPUs have been used to benchmark LLM models:

- A100 which are available on Curnagl, [official documentation](#),
- GH200 which will be available soon on Curnagl, [official documentation](#),
- L40 which will be available soon on Curnagl, [official documentation](#) and [specifications](#).

Here are their specifications

| Characteristics         | A100 | GH200 | L40S |
|-------------------------|------|-------|------|
| Number of nodes at UNIL | 8    | 1     | 8    |
| Memory per node (GB)    | 40   | 80    | 48   |

| Characteristics                            | A100 | GH200 | L40S |
|--|------|-------|------|
| Number of CPU per NUMA node                | 48   | 72    | 8    |
| Memory bandwidth - up to (TB/s)            | 1.9  | 4     | 0.86 |
| FP64 performance (teraFlops)               | 9.7  | 34    | NA   |
| TF64 performance (teraFlops)               | 19.5 | 67    | NA   |
| FP32 performance (teraFlops)               | 19.5 | 67    | 91.6 |
| TF32 performance (teraFlops)               | 156  | 494   | 183  |
| TF32 performance with sparsity (teraFlops) | 312  | 494   | 366  |
| FP16 performance (teraFlops)               | 312  | 990   | 362  |
| INT8 performance (teraFlops)               | 624  | 1.9   | 733  |

Depending on the code you are running, one GPU may better suit your requirements and expectations.

**Note:** These architectures are not powerful enough to train Large Language Models.

**Note:** Our benchmarks aim to determine which GPU types should be provided to researchers. If you require new GPUs for your research, feel free to reach out to us through the Help Desk. In case, you and other researchers agree on the same GPU request, we will do our best to provide new resources that meet your needs.

## Inference latency results

This [chat dataset from GPT3](#) has been used to benchmark models.

In order to guarantee reproducibility of results and be able to perform a comparison between different benchmarks we set the following parameters:

- The maximum number of tokens to generate, is set to
- The temperature, which controls the output randomness, is set to

- The context size, which is the number of tokens the model can process within a single input, is set to `default`. This means the maximum context size of the model (e.g 131072 for Llama3.1)
- Use of GPU exclusively
- All models are loaded in F16 (no quantization)

## Mistral models

### mistral-7B-Instruct-v0.3

| Backend results (Token/seconds) | A100 | GH200 | L40  |
|---------------------------------|------|-------|------|
| vllm                            | 74.1 | -     | -    |
| llama.cpp                       | 53.8 | 138.4 | 42.8 |
| Transformers                    | 30   | 41.3  | 21.6 |
| mistral-inference               | 23.4 | -     | 25   |

### Mixtral-8x7B-v0.1-Instruct

| Backend results (Token/seconds) | A100 | GH200 | L40  |
|---------------------------------|------|-------|------|
| llama.cpp                       | NA   | NA    | 23.4 |
| Transformers                    | NA   | NA    | 8.5  |

## Llama models

### 8B Instruct

| Backend results (Token/seconds) | A100   | GH200   | L40    |
|---------------------------------|--------|---------|--------|
| llama.cpp                       | 62.645 | 100.845 | 43.387 |
| Transformers                    | 31.650 | 43.321  | 21.062 |
| vllm                            | 44.686 | 119.59  | 45.176 |

### 70B Instruct

| Backend results (Token/seconds) | L40 |
|---------------------------------|-----|
|---------------------------------|-----|

|              |        |
|--------------|--------|
| llama.cpp    | 5.029  |
| Transformers | 2.372  |
| vllm         | 30.945 |

# Conclusions

- Mixtral 8x7B and Llama 70B Instruct are composed of several billions of parameters. Therefore the resulting memory consumption for inference can only be supported by multiple GPUs using the same machine or by using a combination of VRAM and RAM host memory. This of course will degrade the performance because we need to transfer data between two types of memory which could be slow. GH200 has a large bus memory which offers a good performance on this types of cases.
- The use of distributed setup adds a lot of latency.
- Transformers backend offers a good trade-off between learning curve and performance.
- Backends offer the possibility to configure a context size. The parameter has no impact on performance (token throughput) but it is correlated to the amount of VRAM consumed. Therefore, if you want to optimize memory consumption you should set the context size to an appropriate value.
- GH200 offers the best inference speed but it could be difficult to set up and install libraries on.
- The results shown here were obtained without any optimization. There are optimization than can be applied like quantization and flash attention.

---

Révision #38

Créé 6 novembre 2024 08:56:43 par Cristian Ruiz

Mis à jour 25 mars 2025 12:39:09 par Cristian Ruiz