

```
In [2]: !pip install --upgrade pip
!pip install -q transformers ipywidgets huggingface_hub
```

Requirement already satisfied: pip in /opt/conda/lib/python3.10/site-packages (23.2.1)

Collecting pip

Obtaining dependency information for pip from <https://files.pythonhosted.org/packages/e7/54/0c1c068542cee73d8863336e974fc881e608d0170f3af15d0c0f28644531/pip-24.1.2-py3-none-any.whl.metadata>

Downloading pip-24.1.2-py3-none-any.whl.metadata (3.6 kB)

Downloading pip-24.1.2-py3-none-any.whl (1.8 MB)

1.8/1.8 MB 20.8 MB/s eta 0:00:0000:0100:01

Installing collected packages: pip

Attempting uninstall: pip

Found existing installation: pip 23.2.1

Uninstalling pip-23.2.1:

Successfully uninstalled pip-23.2.1

Successfully installed pip-24.1.2

WARNING: Running pip as the 'root' user can result in broken permissions and conflicting behaviour with the system package manager. It is recommended to use a virtual environment instead: <https://pip.pypa.io/warnings/venv>

WARNING: Running pip as the 'root' user can result in broken permissions and conflicting behaviour with the system package manager, possibly rendering your system unusable. It is recommended to use a virtual environment instead: <https://pip.pypa.io/warnings/venv>. Use the --root-user-action option if you know what you are doing and want to suppress this warning.

```
In [3]: from huggingface_hub import notebook_login
notebook_login()
```

VBox(children=(HTML(value='<center> <img\\nsrc=https://huggingface.co/front/assets/huggingface_logo-noborder.sv...

```
In [4]: from transformers import AutoTokenizer, AutoModelForCausalLM
```

```
In [5]: tokenizer = AutoTokenizer.from_pretrained("Jacaranda/UlizaLlama3")
model = AutoModelForCausalLM.from_pretrained("Jacaranda/UlizaLlama3")
```

tokenizer_config.json: 0%| | 0.00/50.6k [00:00<?, ?B/s]

tokenizer.json: 0%| | 0.00/9.09M [00:00<?, ?B/s]

special_tokens_map.json: 0%| | 0.00/449 [00:00<?, ?B/s]

config.json: 0%| | 0.00/713 [00:00<?, ?B/s]

model.safetensors.index.json: 0%| | 0.00/23.9k [00:00<?, ?B/s]

Downloading shards: 0%| | 0/4 [00:00<?, ?it/s]

model-00001-of-00004.safetensors: 0%| | 0.00/4.98G [00:00<?, ?B/s]

model-00002-of-00004.safetensors: 0%| | 0.00/5.00G [00:00<?, ?B/s]

model-00003-of-00004.safetensors: 0%| | 0.00/4.92G [00:00<?, ?B/s]

model-00004-of-00004.safetensors: 0%| | 0.00/1.17G [00:00<?, ?B/s]

Loading checkpoint shards: 0%| | 0/4 [00:00<?, ?it/s]

generation_config.json: 0%| | 0.00/172 [00:00<?, ?B/s]

```
In [6]: model
```

```
Out[6]: LlamaForCausalLM(
  (model): LlamaModel(
    (embed_tokens): Embedding(128256, 4096)
    (layers): ModuleList(
      (0-31): 32 x LlamaDecoderLayer(
        (self_attn): LlamaAttention(
          (q_proj): Linear(in_features=4096, out_features=4096, bias=False)
          (k_proj): Linear(in_features=4096, out_features=1024, bias=False)
          (v_proj): Linear(in_features=4096, out_features=1024, bias=False)
          (o_proj): Linear(in_features=4096, out_features=4096, bias=False)
          (rotary_emb): LlamaRotaryEmbedding()
        )
        (mlp): LlamaMLP(
          (gate_proj): Linear(in_features=4096, out_features=14336, bias=False)
          (up_proj): Linear(in_features=4096, out_features=14336, bias=False)
          (down_proj): Linear(in_features=14336, out_features=4096, bias=False)
          (act_fn): SiLU()
        )
        (input_layernorm): LlamaRMSNorm()
        (post_attention_layernorm): LlamaRMSNorm()
      )
    )
    (norm): LlamaRMSNorm()
    (rotary_emb): LlamaRotaryEmbedding()
  )
  (lm_head): Linear(in_features=4096, out_features=128256, bias=False)
)
```

```
In [7]: tokenizer
```

[illegible]

[illegible]

[illegible]

[illegible]

[illegible]

[illegible]

[illegible]

[illegible]

[illegible]

[illegible]

```

special=True),
    128253: AddedToken("<|reserved_special_token_248|>", rstrip=False, lstrip=False, single_word=False, normalized=False,
special=True),
    128254: AddedToken("<|reserved_special_token_249|>", rstrip=False, lstrip=False, single_word=False, normalized=False,
special=True),
    128255: AddedToken("<|reserved_special_token_250|>", rstrip=False, lstrip=False, single_word=False, normalized=False,
special=True),
}

```

In [23]: `!nvidia-smi`

huggingface/tokenizers: The current process just got forked, after parallelism has already been used. Disabling parallelism to a void deadlocks...

To disable this warning, you can either:

- Avoid using `tokenizers` before the fork if possible
- Explicitly set the environment variable TOKENIZERS_PARALLELISM=(true | false)

Sat Jul 27 21:58:57 2024

+-----+									
NVIDIA-SMI 545.23.06				Driver Version: 545.23.06				CUDA Version: 12.3	
+-----+									
GPU		Name		Persistence-M		Bus-Id		Disp.A	
Fan		Temp		Pwr:Usage/Cap		Memory-Usage		Volatile Uncorr. ECC	
		Perf						GPU-Util Compute M.	
								MIG M.	
+=====+									
0		Quadro RTX 8000		On		00000000:0A:00.0 Off		0	
N/A		38C P0		61W / 250W		30803MiB / 46080MiB		0% Default	
								N/A	
+-----+									
+-----+									
Processes:									
GPU		GI		PID		Type		Process name	
		ID						GPU Memory	
		ID						Usage	
+=====+									
+-----+									

In [24]: `import torch`

```

device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
device

```

Out[24]: `device(type='cuda')`

In [26]: `model.to(device)`

```

Out [26]: LlamaForCausalLM(
  (model): LlamaModel(
    (embed_tokens): Embedding(128256, 4096)
    (layers): ModuleList(
      (0-31): 32 x LlamaDecoderLayer(
        (self_attn): LlamaAttention(
          (q_proj): Linear(in_features=4096, out_features=4096, bias=False)
          (k_proj): Linear(in_features=4096, out_features=1024, bias=False)
          (v_proj): Linear(in_features=4096, out_features=1024, bias=False)
          (o_proj): Linear(in_features=4096, out_features=4096, bias=False)
          (rotary_emb): LlamaRotaryEmbedding()
        )
        (mlp): LlamaMLP(
          (gate_proj): Linear(in_features=4096, out_features=14336, bias=False)
          (up_proj): Linear(in_features=4096, out_features=14336, bias=False)
          (down_proj): Linear(in_features=14336, out_features=4096, bias=False)
          (act_fn): SiLU()
        )
        (input_layernorm): LlamaRMSNorm()
        (post_attention_layernorm): LlamaRMSNorm()
      )
    )
    (norm): LlamaRMSNorm()
    (rotary_emb): LlamaRotaryEmbedding()
  )
  (lm_head): Linear(in_features=4096, out_features=128256, bias=False)
)

```

```

In [27]: # Use a pipeline as a high-level helper

```

```

from transformers import pipeline

```

```

pipe = pipeline("question-answering", model = model, tokenizer = tokenizer, device = 0)

```

The model 'LlamaForCausalLM' is not supported for question-answering. Supported models are ['AlbertForQuestionAnswering', 'BartForQuestionAnswering', 'BertForQuestionAnswering', 'BigBirdForQuestionAnswering', 'BigBirdPegasusForQuestionAnswering', 'BloomForQuestionAnswering', 'CamembertForQuestionAnswering', 'CanineForQuestionAnswering', 'ConvBertForQuestionAnswering', 'Data2VecTextForQuestionAnswering', 'DebertaForQuestionAnswering', 'DebertaV2ForQuestionAnswering', 'DistilBertForQuestionAnswering', 'ElectraForQuestionAnswering', 'ErnieForQuestionAnswering', 'ErnieMForQuestionAnswering', 'FalconForQuestionAnswering', 'FlaubertForQuestionAnsweringSimple', 'FNetForQuestionAnswering', 'FunnelForQuestionAnswering', 'GPT2ForQuestionAnswering', 'GPTNeoForQuestionAnswering', 'GPTNeoXForQuestionAnswering', 'GPTJForQuestionAnswering', 'IBertForQuestionAnswering', 'LayoutLMv2ForQuestionAnswering', 'LayoutLMv3ForQuestionAnswering', 'LEDForQuestionAnswering', 'LiltForQuestionAnswering', 'LlamaForQuestionAnswering', 'LongformerForQuestionAnswering', 'LukeForQuestionAnswering', 'LxmertForQuestionAnswering', 'MarkupLMForQuestionAnswering', 'MBartForQuestionAnswering', 'MegaForQuestionAnswering', 'MegatronBertForQuestionAnswering', 'MobileBertForQuestionAnswering', 'MPNetForQuestionAnswering', 'MptForQuestionAnswering', 'MraForQuestionAnswering', 'MT5ForQuestionAnswering', 'MvpForQuestionAnswering', 'NezhaForQuestionAnswering', 'NystromformerForQuestionAnswering', 'OPTForQuestionAnswering', 'QDQBertForQuestionAnswering', 'ReformerForQuestionAnswering', 'RemBertForQuestionAnswering', 'RobertaForQuestionAnswering', 'RobertaPreLayerNormForQuestionAnswering', 'RoCBertForQuestionAnswering', 'RoFormerForQuestionAnswering', 'SplinterForQuestionAnswering', 'SqueezeBertForQuestionAnswering', 'T5ForQuestionAnswering', 'UMT5ForQuestionAnswering', 'XLNetForQuestionAnsweringSimple', 'XLMRobertaForQuestionAnswering', 'XLMRobertaXLForQuestionAnswering', 'XLNetForQuestionAnsweringSimple', 'XmodForQuestionAnswering', 'YosoForQuestionAnswering'].

```
In [28]: question = "What is the capital of France?"
```

```
In [29]: # Preprocess the question
inputs = tokenizer(question, return_tensors="pt")
```

```
In [30]: inputs['attention_mask']
```

```
Out[30]: tensor([[1, 1, 1, 1, 1, 1, 1, 1]])
```

```
In [32]: # Move the inputs to the specified device (GPU or CPU)
inputs = {k: v.to(device) for k, v in inputs.items()}
```

```
In [33]: # Generate an answer
outputs = model.generate(inputs["input_ids"], attention_mask=inputs["attention_mask"])

# roughly 6 minutes on RTX 8000 X1
```

```
In [34]: # Convert the answer to text
answer = tokenizer.decode(outputs[0].cpu(), skip_special_tokens=True)

print(answer)
```


[illegible]

[illegible]

[illegible]

In []: