

Token Generation Pipeline for LLM

Use case with HuggingFace

Objectifs

Due to the revolution of LLM model in the world of AI, companies want to integrate its performances into their services.

However the use of the API can be expensive and the models are too large to fit on regular devices. The solution will be to split parameters over to disk and ram.

Documentations

The solution is divided on several parts and parameters and we use some high level packages developed by HuggingFace based on the following documentations:

1. <https://huggingface.co/blog/accelerate-large-models>
2. https://huggingface.co/docs/accelerate/usage_guides/big_modeling
3. https://huggingface.co/docs/transformers/main/main_classes/quantization

Solution

Many aspects of the model have to be taken with consideration like the parameters format and its architecture.

However, when you load a model directly with HuggingFace package, you load the architecture and the parameters.

So, the main idea is to start with an empty model and make sure on which device to store each layer (submodule).

1. We define the checkpoint of the model we want to load.
2. We load the empty model
3. We define the device mapping
4. We configure a quantization object
5. We load the model from HuggingFace set up with the previous steps and an offload folder.

Example

We load the model Bloom-3B, with 3 billions of parameters. It needs 24 Go in the RAM to be able to be loaded. However, the RAM of the T4 colab gpu is 15 Go.

The solution organize the memory ingestion between the CPU RAM, GPU RAM and the drive.

```
1  # define model checkpoint
2  checkpoint = "bigscience/bloom-3b"
3
4  # define no splitting block name for device mapping
5  no_split_block = 'BloomBlock'
6
7  # load empty model
8  config = AutoConfig.from_pretrained(checkpoint)
9  with init_empty_weights():
10     model = AutoModelForCausalLM.from_config(config)
11
12  # device block mapping
13  device_map = infer_auto_device_map(model,
14                                     no_split_module_classes=[no_split_block],
15                                     dtype='float16')
16
17  # set up a quantization
18  quantization_config = BitsAndBytesConfig(llm_int8_enable_fp32_cpu_offload=True,
19                                           llm_int8_threshold=6.0,
20                                           llm_int8_skip_modules=["lm_head"])
21
22  # load model from sharded files
23  model = AutoModelForCausalLM.from_pretrained(checkpoint,
24                                               device_map=device_map,
25                                               offload_folder="offload",
26                                               offload_state_dict=True,
27                                               load_in_8bit=True,
28                                               quantization_config=quantization_config,
29                                               )
30
31  # load tokenizer
32  tokenizer = AutoTokenizer.from_pretrained(checkpoint)
33
34  # generate inputs
35  prompt = "The quick brown fox"
36  max_length = 50
37  inputs = tokenizer(prompt, return_tensors="pt").to('cuda')
38
39  # get outputs
40  outputs = model.generate(inputs["input_ids"], max_length=max_length)
41
42  # decode outputs
43  decoded_outputs = tokenizer.batch_decode(outputs, skip_special_tokens=True)
44
45  # print outputs
46  print(decoded_outputs)
```

Conclusion

This work was done on a very short amount of time (1 week) and I needed to explore very new technologies that never stops to be improved.

With the GPU device offered by Google Colab the pipeline works for “small” LLM models - with more than 6B of parameters the colab session crashes.

The pipeline depends a lot on the parameters format and can lose a part of its performance.