

```
1 !pip install git+https://github.com/shumingma/transformers.git
```

```
2
```

```
Collecting git+https://github.com/shumingma/transformers.git
  Cloning https://github.com/shumingma/transformers.git to /tmp/pip-req-build-gaq_vty7
  Running command git clone --filter=blob:none --quiet https://github.com/shumingma/transformers.git /tmp/pip-req-build-gaq_vty7
  Resolved https://github.com/shumingma/transformers.git to commit 4a01efe84d0120dc2545ff2de445082400d87407
  Installing build dependencies ... done
  Getting requirements to build wheel ... done
  Preparing metadata (pyproject.toml) ... done
Requirement already satisfied: filelock in /usr/local/lib/python3.11/dist-packages (from transformers==4.52.0.dev0) (3.18.0)
Requirement already satisfied: huggingface-hub<1.0,>=0.30.0 in /usr/local/lib/python3.11/dist-packages (from transformers==4.52.0.dev0) (2.0.2)
Requirement already satisfied: numpy>=1.17 in /usr/local/lib/python3.11/dist-packages (from transformers==4.52.0.dev0) (2.0.2)
Requirement already satisfied: packaging>=20.0 in /usr/local/lib/python3.11/dist-packages (from transformers==4.52.0.dev0) (24.2)
Requirement already satisfied: pyyaml>=5.1 in /usr/local/lib/python3.11/dist-packages (from transformers==4.52.0.dev0) (6.0.2)
Requirement already satisfied: regex!=2019.12.17 in /usr/local/lib/python3.11/dist-packages (from transformers==4.52.0.dev0) (2024.11.6)
Requirement already satisfied: requests in /usr/local/lib/python3.11/dist-packages (from transformers==4.52.0.dev0) (2.32.3)
Requirement already satisfied: tokenizers<0.22,>=0.21 in /usr/local/lib/python3.11/dist-packages (from transformers==4.52.0.dev0) (0.21.
Requirement already satisfied: safetensors>=0.4.3 in /usr/local/lib/python3.11/dist-packages (from transformers==4.52.0.dev0) (0.5.3)
Requirement already satisfied: tqdm>=4.27 in /usr/local/lib/python3.11/dist-packages (from transformers==4.52.0.dev0) (4.67.1)
Requirement already satisfied: fsspec>=2023.5.0 in /usr/local/lib/python3.11/dist-packages (from huggingface-hub<1.0,>=0.30.0->transformers==4.52.0.dev0) (2024.10.1)
Requirement already satisfied: typing-extensions>=3.7.4.3 in /usr/local/lib/python3.11/dist-packages (from huggingface-hub<1.0,>=0.30.0->transformers==4.52.0.dev0) (4.12.2)
Requirement already satisfied: charset-normalizer<4,>=2 in /usr/local/lib/python3.11/dist-packages (from requests->transformers==4.52.0.dev0) (3.10.0)
Requirement already satisfied: idna<4,>=2.5 in /usr/local/lib/python3.11/dist-packages (from requests->transformers==4.52.0.dev0) (3.10.0)
Requirement already satisfied: urllib3<3,>=1.21.1 in /usr/local/lib/python3.11/dist-packages (from requests->transformers==4.52.0.dev0) (2.3.0)
Requirement already satisfied: certifi>=2017.4.17 in /usr/local/lib/python3.11/dist-packages (from requests->transformers==4.52.0.dev0) (2025.1.1)
Building wheels for collected packages: transformers
  Building wheel for transformers (pyproject.toml) ... done
  Created wheel for transformers: filename=transformers-4.52.0.dev0-py3-none-any.whl size=11413940 sha256=1111ae15305127dd855dc5b8049916
  Stored in directory: /tmp/pip-ephem-wheel-cache-fxk3dmlj/wheels/2c/2a/7c/3be0c30fb51a7becc4bcb536739ae9ed9cc7e633ffbfaf63b
Successfully built transformers
Installing collected packages: transformers
  Attempting uninstall: transformers
    Found existing installation: transformers 4.51.3
    Uninstalling transformers-4.51.3:
      Successfully uninstalled transformers-4.51.3
  Successfully installed transformers-4.52.0.dev0
```

```
1 !git clone https://huggingface.co/microsoft/bitnet-b1.58-2B-4T
```

```
Cloning into 'bitnet-b1.58-2B-4T'...
remote: Enumerating objects: 46, done.
remote: Counting objects: 100% (45/45), done.
remote: Compressing objects: 100% (43/43), done.
remote: Total 46 (delta 18), reused 0 (delta 0), pack-reused 1 (from 1)
Unpacking objects: 100% (46/46), 2.23 MiB | 5.10 MiB/s, done.
```

```
1 import torch
2 from transformers import AutoModelForCausalLM, AutoTokenizer
3
4 # Specify the model path
5 model_path = "/content/bitnet-b1.58-2B-4T"
6
7 # Load tokenizer (still from the model ID)
8 tokenizer = AutoTokenizer.from_pretrained("microsoft/bitnet-b1.58-2B-4T")
9
10 # Load model from the specified path
11 model = AutoModelForCausalLM.from_pretrained(
12     model_path,
13     torch_dtype=torch.bfloat16
14 )
15
16 # ... (rest of your code remains the same)
```

```
tokenizer_config.json: 100% 50.8k/50.8k [00:00<00:00, 957kB/s]
tokenizer.json: 100% 9.09M/9.09M [00:00<00:00, 25.9MB/s]
special_tokens_map.json: 100% 73.0/73.0 [00:00<00:00, 2.07kB/s]
You don't have a GPU available to load the model, the inference will be slow because of weight unpacking
No CUDA runtime is found. using CUDA_HOME='/usr/local/cuda'
```

```
1 # Apply the chat template
2 messages = [
3     {"role": "system", "content": "You are a helpful AI assistant."},
4     {"role": "user", "content": "How are you?"},
5 ]
```

```

5 ]
6 prompt = tokenizer.apply_chat_template(messages, tokenize=False, add_generation_prompt=True)
7 chat_input = tokenizer(prompt, return_tensors="pt").to(model.device)
8
9 # Generate response
10 chat_outputs = model.generate(**chat_input, max_new_tokens=11)
11 response = tokenizer.decode(chat_outputs[0][chat_input['input_ids'].shape[-1]:], skip_special_tokens=True) # Decode only the response part
12 print("\nAssistant Response:", response)

```

Setting `pad_token_id` to `eos_token_id`:128001 for open-end generation.

Assistant Response: I'm just a computer program, so I don't

1 Start coding or [generate](#) with AI.

1 Start coding or [generate](#) with AI.

1 Start coding or [generate](#) with AI.

✓ /content/bitnet-b1.58-2B-4T/configuration_bitnet.py

```

1 #/content/bitnet-b1.58-2B-4T/configuration_bitnet.py
2
3
4
5 # coding=utf-8
6 # Copyright 2022 EleutherAI and the HuggingFace Inc. team. All rights reserved.
7 #
8 # This code is based on EleutherAI's GPT-NeoX library and the GPT-NeoX
9 # and OPT implementations in this library. It has been modified from its
10 # original forms to accommodate minor architectural differences compared
11 # to GPT-NeoX and OPT used by the Meta AI team that trained the model.
12 #
13 # Licensed under the Apache License, Version 2.0 (the "License");
14 # you may not use this file except in compliance with the License.
15 # You may obtain a copy of the License at
16 #
17 #     http://www.apache.org/licenses/LICENSE-2.0
18 #
19 # Unless required by applicable law or agreed to in writing, software
20 # distributed under the License is distributed on an "AS IS" BASIS,
21 # WITHOUT WARRANTIES OR CONDITIONS OF ANY KIND, either express or implied.
22 # See the License for the specific language governing permissions and
23 # limitations under the License.
24 """ LLaMA model configuration"""
25
26 from transformers.configuration_utils import PretrainedConfig
27 from transformers.utils import logging
28
29
30 logger = logging.get_logger(__name__)
31
32 LLAMA_PRETRAINED_CONFIG_ARCHIVE_MAP = {}
33
34
35 class BitnetConfig(PretrainedConfig):
36     r"""
37     This is the configuration class to store the configuration of a [`BitnetModel`]. It is used to instantiate an LLaMA
38     model according to the specified arguments, defining the model architecture. Instantiating a configuration with the
39     defaults will yield a similar configuration to that of the LLaMA-7B.
40     Configuration objects inherit from [`PretrainedConfig`] and can be used to control the model outputs. Read the
41     documentation from [`PretrainedConfig`] for more information.
42
43     Args:
44         vocab_size (`int`, *optional*, defaults to 32000):
45             Vocabulary size of the LLaMA model. Defines the number of different tokens that can be represented by the
46             `inputs_ids` passed when calling [`BitnetModel`]
47         hidden_size (`int`, *optional*, defaults to 4096):
48             Dimension of the hidden representations.
49         intermediate_size (`int`, *optional*, defaults to 11008):
50             Dimension of the MLP representations.
51         num_hidden_layers (`int`, *optional*, defaults to 32):
52             Number of hidden layers in the Transformer decoder.
53         num_attention_heads (`int`, *optional*, defaults to 32):

```

```

53         Number of attention heads for each attention layer in the Transformer decoder.
54     num_key_value_heads (`int`, *optional*):
55         This is the number of key_value heads that should be used to implement Grouped Query Attention. If
56         `num_key_value_heads=num_attention_heads`, the model will use Multi Head Attention (MHA), if
57         `num_key_value_heads=1` the model will use Multi Query Attention (MQA) otherwise GQA is used. When
58         converting a multi-head checkpoint to a GQA checkpoint, each group key and value head should be constructed
59         by meanpooling all the original heads within that group. For more details checkout [this
60         paper](https://arxiv.org/pdf/2305.13245.pdf). If it is not specified, will default to
61         `num_attention_heads`.
62     hidden_act (`str` or `function`, *optional*, defaults to `"silu"`):
63         The non-linear activation function (function or string) in the decoder.
64     max_position_embeddings (`int`, *optional*, defaults to 2048):
65         The maximum sequence length that this model might ever be used with. Bitnet 1 supports up to 2048 tokens,
66         Bitnet 2 up to 4096, CodeBitnet up to 16384.
67     initializer_range (`float`, *optional*, defaults to 0.02):
68         The standard deviation of the truncated_normal_initializer for initializing all weight matrices.
69     rms_norm_eps (`float`, *optional*, defaults to 1e-06):
70         The epsilon used by the rms normalization layers.
71     use_cache (`bool`, *optional*, defaults to `True`):
72         Whether or not the model should return the last key/values attentions (not used by all models). Only
73         relevant if `config.is_decoder=True`.
74     pad_token_id (`int`, *optional*):
75         Padding token id.
76     bos_token_id (`int`, *optional*, defaults to 1):
77         Beginning of stream token id.
78     eos_token_id (`int`, *optional*, defaults to 2):
79         End of stream token id.
80     pretraining_tp (`int`, *optional*, defaults to 1):
81         Experimental feature. Tensor parallelism rank used during pretraining. Please refer to [this
82         document](https://huggingface.co/docs/transformers/main/perf_train_gpu_many#tensor-parallelism) to understand more about it
83         necessary to ensure exact reproducibility of the pretraining results. Please refer to [this
84         issue](https://github.com/pytorch/pytorch/issues/76232).
85     tie_word_embeddings (`bool`, *optional*, defaults to `False`):
86         Whether to tie weight embeddings
87     rope_theta (`float`, *optional*, defaults to 10000.0):
88         The base period of the RoPE embeddings.
89     rope_scaling (`Dict`, *optional*):
90         Dictionary containing the scaling configuration for the RoPE embeddings. Currently supports two scaling
91         strategies: linear and dynamic. Their scaling factor must be a float greater than 1. The expected format is
92         `{"type": strategy name, "factor": scaling factor}`. When using this flag, don't update
93         `max_position_embeddings` to the expected new maximum. See the following thread for more information on how
94         these scaling strategies behave:
95         https://www.reddit.com/r/LocalLLaMA/comments/14mrgpr/dynamically_scaled_rope_further_increases/. This is an
96         experimental feature, subject to breaking API changes in future versions.
97     attention_bias (`bool`, defaults to `False`, *optional*, defaults to `False`):
98         Whether to use a bias in the query, key, value and output projection layers during self-attention.
99     attention_dropout (`float`, *optional*, defaults to 0.0):
100         The dropout ratio for the attention probabilities.
101
102     """
103     python
104     >>> from transformers import BitnetModel, BitnetConfig
105     >>> # Initializing a LLaMA llama-7b style configuration
106     >>> configuration = BitnetConfig()
107     >>> # Initializing a model from the llama-7b style configuration
108     >>> model = BitnetModel(configuration)
109     >>> # Accessing the model configuration
110     >>> configuration = model.config
111     >>> """
112
113     model_type = "llama"
114     keys_to_ignore_at_inference = ["past_key_values"]
115
116     def __init__(
117         self,
118         vocab_size=32000,
119         hidden_size=4096,
120         intermediate_size=11008,
121         num_hidden_layers=32,
122         num_attention_heads=32,
123         num_key_value_heads=None,
124         hidden_act="silu",
125         max_position_embeddings=2048,
126         initializer_range=0.02,
127         rms_norm_eps=1e-6,
128         use_cache=True,
129         pad_token_id=None,
130         bos_token_id=1,
131         eos_token_id=2,

```

```

130     pretraining_tp=1,
131     tie_word_embeddings=False,
132     rope_theta=10000.0,
133     rope_scaling=None,
134     attention_bias=False,
135     attention_dropout=0.0,
136     weight_bits=1,
137     input_bits=8,
138     **kwargs,
139 ):
140     self.vocab_size = vocab_size
141     self.max_position_embeddings = max_position_embeddings
142     self.hidden_size = hidden_size
143     self.intermediate_size = intermediate_size
144     self.num_hidden_layers = num_hidden_layers
145     self.num_attention_heads = num_attention_heads
146
147     # for backward compatibility
148     if num_key_value_heads is None:
149         num_key_value_heads = num_attention_heads
150
151     self.num_key_value_heads = num_key_value_heads
152     self.hidden_act = hidden_act
153     self.initializer_range = initializer_range
154     self.rms_norm_eps = rms_norm_eps
155     self.pretraining_tp = pretraining_tp
156     self.use_cache = use_cache
157     self.rope_theta = rope_theta
158     self.rope_scaling = rope_scaling
159     self._rope_scaling_validation()
160     self.attention_bias = attention_bias
161     self.attention_dropout = attention_dropout
162     self.weight_bits = weight_bits
163     self.input_bits = input_bits
164
165     super().__init__(
166         pad_token_id=pad_token_id,
167         bos_token_id=bos_token_id,
168         eos_token_id=eos_token_id,
169         tie_word_embeddings=tie_word_embeddings,
170         **kwargs,
171     )
172
173     def _rope_scaling_validation(self):
174         """
175         Validate the `rope_scaling` configuration.
176         """
177         if self.rope_scaling is None:
178             return
179
180         if not isinstance(self.rope_scaling, dict) or len(self.rope_scaling) != 2:
181             raise ValueError(
182                 "`rope_scaling` must be a dictionary with with two fields, `type` and `factor`, "
183                 f"got {self.rope_scaling}"
184             )
185         rope_scaling_type = self.rope_scaling.get("type", None)
186         rope_scaling_factor = self.rope_scaling.get("factor", None)
187         if rope_scaling_type is None or rope_scaling_type not in ["linear", "dynamic"]:
188             raise ValueError(
189                 f"`rope_scaling`'s type field must be one of ['linear', 'dynamic'], got {rope_scaling_type}"
190             )
191         if rope_scaling_factor is None or not isinstance(rope_scaling_factor, float) or rope_scaling_factor <= 1.0:
192             raise ValueError(f"`rope_scaling`'s factor field must be a float > 1, got {rope_scaling_factor}")

```

1 Start coding or [generate](#) with AI.

1 Start coding or [generate](#) with AI.

1 Start coding or [generate](#) with AI.

```

1 from huggingface_hub import HfApi, upload_folder
2
3 # إعدادات المستخدم
4 token = "XXXXXXXXXXXXXXXXXXXXXXXXXXXX"

```

```

5 repo_id = "rakmik/bitnetrun"
6 folder_path = "/content/bitnet-b1.58-2B-4T" # غير هذا إلى مسار مجلدك المحلي ←
7
8 # رفع المجلد
9 upload_folder(
10     repo_id=repo_id,
11     folder_path=folder_path,
12     path_in_repo="", # يمكنك تحديد مجلد داخل الريبو إذا رغبت
13     repo_type="model", # "dataset" أو
14     token=token,
15 )
16 print("✅ تم رفع المجلد بنجاح إلى ", repo_id)
17

```

model.safetensors: 100% 1.18G/1.18G [00:29<00:00, 44.6MB/s]

✅ تم رفع المجلد بنجاح إلى rakmik/bitnetrun

1 Start coding or [generate](#) with AI.

1 Start coding or [generate](#) with AI.

1 Start coding or [generate](#) with AI.

```

1 import torch
2 from transformers import AutoModelForCausalLM, AutoTokenizer
3
4 model_id = "rakmik/bitnetrun"
5
6 # Load tokenizer and model
7 tokenizer = AutoTokenizer.from_pretrained(model_id)
8 model = AutoModelForCausalLM.from_pretrained(
9     model_id,
10     torch_dtype=torch.bfloat16
11 )
12
13 # Apply the chat template
14 messages = [
15     {"role": "system", "content": "You are a helpful AI assistant."},
16     {"role": "user", "content": "How are you?"},
17 ]
18 prompt = tokenizer.apply_chat_template(messages, tokenize=False, add_generation_prompt=True)
19 chat_input = tokenizer(prompt, return_tensors="pt").to(model.device)
20
21 # Generate response
22 chat_outputs = model.generate(**chat_input, max_new_tokens=11)
23 response = tokenizer.decode(chat_outputs[0][chat_input['input_ids'].shape[-1]:], skip_special_tokens=True) # Decode only the response part
24 print("\nAssistant Response:", response)
25

```

tokenizer_config.json: 100% 50.8k/50.8k [00:00<00:00, 1.03MB/s]

tokenizer.json: 100% 9.09M/9.09M [00:00<00:00, 24.3MB/s]

special_tokens_map.json: 100% 73.0/73.0 [00:00<00:00, 6.08kB/s]

config.json: 100% 803/803 [00:00<00:00, 68.3kB/s]

model.safetensors: 92% 1.08G/1.18G [00:30<00:00, 105MB/s]

Error while downloading from <https://cdn-lfs-us-1.hf.co/repos/a3/7b/a37bcf8eca6b5df2fea1359acb81db8b715d3c720d5857750b05c83c7892ee12/814>
Trying to resume download...

WARNING:huggingface_hub.file_download:Error while downloading from <https://cdn-lfs-us-1.hf.co/repos/a3/7b/a37bcf8eca6b5df2fea1359acb81db8b715d3c720d5857750b05c83c7892ee12/814>
Trying to resume download...

model.safetensors: 100% 1.18G/1.18G [00:01<00:00, 109MB/s]

generation_config.json: 100% 199/199 [00:00<00:00, 11.3kB/s]

Setting `pad_token_id` to `eos_token_id`:128001 for open-end generation.

Assistant Response: As an AI, I don't have feelings, but

1 Start coding or [generate](#) with AI.

1 Start coding or [generate](#) with AI.

1 Start coding or [generate](#) with AI.

```

1 # Apply the chat template
2 messages = [
3     {"role": "system", "content": "You are a helpful AI assistant."},
4     {"role": "user", "content": "how is python?"},
5 ]
6 prompt = tokenizer.apply_chat_template(messages, tokenize=False, add_generation_prompt=True)
7 chat_input = tokenizer(prompt, return_tensors="pt").to(model.device)
8
9 # Generate response
10 chat_outputs = model.generate(**chat_input, max_new_tokens=55)
11 response = tokenizer.decode(chat_outputs[0][chat_input['input_ids'].shape[-1]:], skip_special_tokens=True) # Decode only the response part
12 print("\nAssistant Response:", response)

```

⚡ Setting `pad_token_id` to `eos_token_id`:128001 for open-end generation.

Assistant Response: Python is a high-level, interpreted programming language known for its simplicity and readability. It is widely used

Here are some of its key features:

1. ****Easy to Learn****: Python has a



1 Start coding or [generate](#) with AI.