

# Transformers による GGUF モデルの高速化

インテル® プラットフォームでのパフォーマンスと  
メモリー使用量の改善

Bo Dong、Jun Lin、Zhentao Yu、Zhenzhong Xu、Yu Luo、Hanwen Chang、Haihao Shen  
インテル コーポレーション

GGUF (GPT-Generated Unified Format、GPT によって生成された統一形式) は、ファイル内のテンソルとメタデータを素早く検査できる新しいバイナリー形式です (図 1)。これは言語モデルファイル形式の大きな飛躍であり、GPT などの大規模言語モデル (LLM) の保存と処理の効率を最適化します。PyTorch\* モデルを GGUF 形式に変換するのは簡単です。

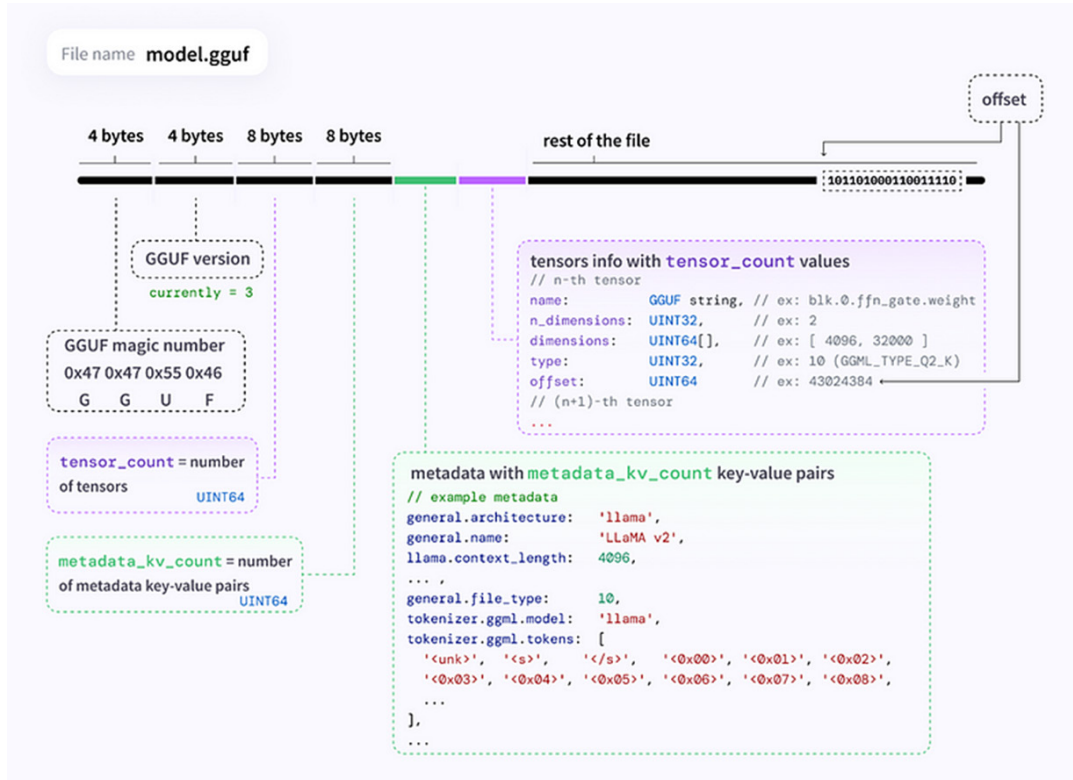


図 1. GGUF 形式 (出典 (英語))

Hugging Face\* Transformers は最近、[Transformers PR](#) (英語)で GGUF をサポートしました。Transformers は、PyTorch\* を使用して推論を実行する前に、GGUF モデルを FP32 に逆量子化します。使い方は簡単で、[図 2](#) に示すように、`from_pretrained` で `gguf_file` パラメーターを指定するだけです。

```
from transformers import AutoTokenizer, TextStreamer, AutoModelForCausalLM

model_name = "TheBloke/TinyLlama-1.1B-Chat-v1.0-GGUF"
gguf_file = "tinylama-1.1b-chat-v1.0.Q4_0.gguf"
model = AutoModelForCausalLM.from_pretrained(model_name, gguf_file = gguf_file)
```

図 2. Hugging Face\* Transformers での GGUF モデルの使用

[Transformers 向けインテル® エクステンション](#) (英語) は低ビットの LLM 推論を高速化します。Hugging Face\* Transformers を拡張し、インテル® プラットフォーム上でパフォーマンスを向上します。幅広いモデルで GGUF 推論をサポートしており、使い方も簡単です (図 3)。

```

from intel_extension_for_transformers.transformers import AutoModelForCausalLM
from transformers import AutoTokenizer, TextStreamer

model_name = "TheBloke/TinyLlama-1.1B-Chat-v1.0-GGUF"
gguf_file = "tinylama-1.1b-chat-v1.0.Q4_0.gguf"
model = AutoModelForCausalLM.from_pretrained(model_name, gguf_file = gguf_file)
    
```

図 3. [Transformers 向けインテル® エクステンション](#) (英語) での GGUF モデルの使用

現在、Transformers は Llama\* や Mistral\* など、50 を超える一般的な LLM をサポートしています (表 1)。

モデル	LLM	中国語モデル	コーディング・モデル	音声モデル
	<a href="#">Meta-Llama-3-8B-Instruct</a> , <a href="#">TinyLlama-1.1B</a> , <a href="#">LLaMA2-tB</a> , <a href="#">LLaMA2-13B</a> , <a href="#">LLaMA2-70B</a> , <a href="#">LLaMA-7B</a> , <a href="#">LLaMA-13B</a> , <a href="#">Solar-10.7B</a> , <a href="#">Mistral-7B</a> , <a href="#">Mistral-7B-Instruct-v0.2</a> , <a href="#">Mixtral-8x7B</a> , <a href="#">GPT-J-6B</a> , <a href="#">GPT-NeoX-20B</a> , <a href="#">Dolly-v2-3B</a> , <a href="#">MPT-7B</a> , <a href="#">MPT-30B</a> , <a href="#">Falcon-7B</a> , <a href="#">Falcon-40B</a> , <a href="#">BLOOM-7B</a> , <a href="#">OPT-125m</a> , <a href="#">OPT-1.3B</a> , <a href="#">OPT-13B</a> , <a href="#">phi-2</a> , <a href="#">phi-1.5</a> , <a href="#">phi-1</a> , <a href="#">phi-3-128k</a> , <a href="#">phi-3-48k</a> , <a href="#">StableLM-2-1.6B</a> , <a href="#">StableLM-3B</a> , <a href="#">StableLM-2-12B</a> , <a href="#">gemma-2b-it</a> , <a href="#">gemma-7b</a> , <a href="#">Neural-Chat-7B-v3-1</a> , <a href="#">Neural-Chat-7B-v3-2</a> ,	<a href="#">Qwen-7B</a> , <a href="#">Qwen-14B</a> , <a href="#">Qwen1.5-7B</a> , <a href="#">Qwen1.5-0.5B</a> , <a href="#">ChatGLM-6B</a> , <a href="#">ChatGLM2-6B</a> , <a href="#">ChatGLM3-6B</a> , <a href="#">Baichuan-13B-Chat</a> , <a href="#">Baichuan2-13B-Chat</a> , <a href="#">Baichuan2-7B-Chat</a>	<a href="#">Magicoder-6.7B</a> , <a href="#">StarCoder-1B</a> , <a href="#">StarCoder-3B</a> , <a href="#">StarCoder-15.5B</a> , <a href="#">CodeLlama-7b</a>	<a href="#">Whisper-tiny</a> , <a href="#">Whisper-base</a> , <a href="#">Whisper-small</a> , <a href="#">Whisper-medium</a> , <a href="#">Whisper-large</a>

表 1. サポートされる LLM



セットアップは簡単です。最初に、Neural Speed をインストールします。

```
git clone https://github.com/intel/neural-speed.git
cd neural_speed
python setup.py install
```

次に、[Transformers 向けインテル® エクステンション](#) (英語) をインストールします。

```
git clone https://github.com/intel/intel-extension-for-transformers.git
cd intel_extension_for_transformers
python setup.py install
```

最後に、Transformers をインストールします。

```
pip install transformers
```

図 4 のコードを使用して簡単にテキストを生成できます。

```
from transformers import AutoTokenizer, TextStreamer
from intel_extension_for_transformers.transformers import AutoModelForCausalLM

# Specify the GGUF repo on the Huggingface
model_name = "TheBloke/TinyLlama-1.1B-Chat-v1.0-GGUF"
# Download the the specific gguf model file from the above repo
gguf_file = "tinyllama-1.1b-chat-v1.0.Q4_0.gguf"
# make sure you are granted to access this model on the Huggingface.
tokenizer_name = "TinyLlama/TinyLlama-1.1B-Chat-v1.0"
prompt = "Once upon a time, there existed a little girl,"
tokenizer = AutoTokenizer.from_pretrained(tokenizer_name, trust_remote_code=True)
inputs = tokenizer(prompt, return_tensors="pt").input_ids
streamer = TextStreamer(tokenizer)
model = AutoModelForCausalLM.from_pretrained(model_name, gguf_file = gguf_file)
outputs = model.generate(inputs, streamer=streamer, max_new_tokens=300)
```

図 4. テキストを生成するサンプルコード

インテル® Xeon® Platinum プロセッサ / インテル® Core™ Ultra 5 プロセッサ 125H ベースのシステム (最大 4400MHz、最小 400MHz、合計メモリー 32GB (8 x 4GB LPDDR5、[7467 MT/s])、Ubuntu\* 13.2.0-9ubuntu1) で Transformers のパフォーマンス比較を実施しました。メモリーが小さいため、比較は [TinyLlama-1.1B-Chat](#) (英語) に限定されました。Transformers 向けインテル® エクステンションは、Transformers と比較して優れたレイテンシーとメモリー使用量を提供しました (表 2)。

TinyLlama in=32, out=32	最初のトークンの レイテンシー	次のトークンの レイテンシー	メモリー使用量	Lambda 精度
Transformers 向け インテル® エクステンション	275ms	12.15ms	641MB	59.07
Hugging Face* Transformers	302ms	88.7ms	4416MB	58.9

表 2. パフォーマンス、メモリー使用量、精度の比較

Transformers 向けインテル® エクステンションは、わずかなコード変更で、メモリー制約のあるシステムでトークンのレイテンシーを最小限に抑える優れた方法です。カーネルは 1 ~ 8 ビットをサポートします。StreamingLLM や Tensor Parallelism などのより高度な機能もあるため、[Transformers 向けインテル® エクステンション](#) (英語) を確認することを推奨します。