

Llama.cpp ausprobieren

Beim Versuch, `llama.cpp` mit einem Modell auszuführen, könnten Sie auf einen Fehler wie diesen stoßen:

```
(py311) lzwjava@Zhiweis-MacBook-Air llama.cpp % ./main -m models/7B/Phi-3-mini-4k-instruct-q4.gguf
main: build = 964 (f3c3b4b)
main: seed  = 1737736417
llama.cpp: loading model from models/7B/Phi-3-mini-4k-instruct-q4.gguf
error loading model: unknown (magic, version) combination: 46554747, 00000003; is this really a GGML file?
llama_load_model_from_file: failed to load model
llama_init_from_gpt_params: error: failed to load model 'models/7B/Phi-3-mini-4k-instruct-q4.gguf'
main: error: unable to load model
```

Dieser Fehler deutet in der Regel auf ein Problem mit der Installation von `llama.cpp` oder der Modell-Datei selbst hin.

Eine gängige Lösung besteht darin, `llama.cpp` mit Homebrew zu installieren:

```
brew install llama.cpp
```

Dies stellt sicher, dass Sie eine kompatible Version der Bibliothek verwenden.

Hier sind einige nützliche Ressourcen:

- Hugging Face GGML Models
- Llama.cpp GitHub Repository
- ggml GitHub Repository
- Ollama
- Ollamac