

Exercise 01

Constrain the output with a custom grammar

1. Overview

GBNF (GGML BNF) is a format for defining formal grammars to constrain model outputs in llama.cpp. For example, you can use it to force the model to generate valid JSON, or speak only in emojis. GBNF grammars are supported in various ways in examples/main and examples/server.

2. Requirements

a) Please read the instructions at the link:

<https://github.com/ggerganov/llama.cpp/blob/master/grammars/README.md>

b) Understand the content in the link and summarize how to constrain the output with a custom grammar.

c) Submit the summarization in one A4 page (pdf file).