

様々な大規模言語モデル (LLM) を実行・評価するための環境構築

千葉啓和 (ライフサイエンス統合データベースセンター)

Methods

クラウド環境、およびオンプレミス環境において様々なLLMを一度に実行し、結果を比較できる環境を構築した。フロントエンドはJavaScript、バックエンドはPython (Flask) で実装した。プログラムからAPIにアクセスすることも可能になっている。

Microsoft Azure OpenAI

- gpt-4o
- gpt-4o-mini
- gpt-4
- gpt-35-turbo



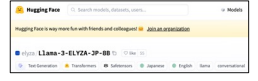
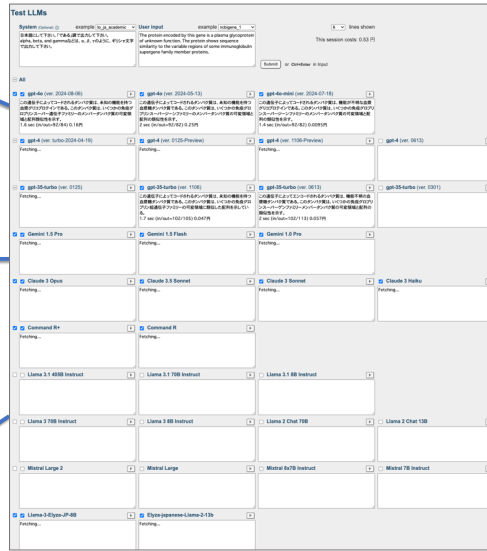
Google Gemini

- Gemini 1.5 Pro
- Gemini 1.5 Flash
- Gemini 1.0 Pro



Amazon Bedrock

- Claude
- Command R
- Llama
- Mistral



Hugging Face

オープンソース
モデル

Elyza-japanese
etc.



オンプレミス
GPUサーバー
(32GB VRAM x 8枚)

ここ1~2年で、実に多くの種類の大規模言語モデル (LLM) が発表された。LLMを活用することによって、生命科学分野におけるデータ処理が加速されると考えられる。しかしながら、実際に各種のLLMを利用してみると、モデルやバージョンによって出力が異なり、どのLLMを用いればよいか分からないことも多い。そこで今回、様々なLLMを一度に実行し、出力を比較できるような環境の構築を行った。LLMの実行環境は、クラウドとGPUサーバーを併用した。より具体的には、OpenAIのモデルに関してはMicrosoft Azureを利用した。Anthropic、Cohere、Meta、Mistral AIのモデルに関しては、Amazon Bedrockを利用した。Geminiに関しては、Googleにアクセスしている。それ以外のモデルに関しては、Hugging Faceからモデルをダウンロードして、GPUサーバー上で実行した。現状では、33種類のモデル (バージョンの違いも含む) を試した結果を一度に表示することができている。ここでは、実行結果とその評価についても共有し議論したい。

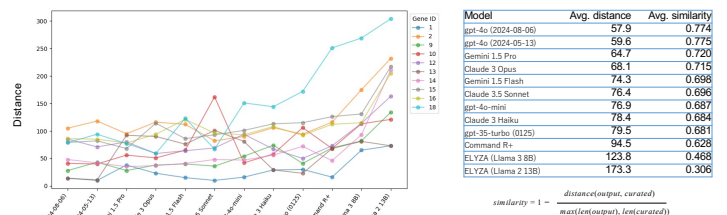
Results

ヒトのNCBI GeneのSummaryを、各LLMを用いて和訳した。一方で、10個の遺伝子について、キュレーションされた和訳も作成し、各モデルの出力との違いを定量的に評価した。差分の可視化にはJavaScriptのjsdiffを用い、距離の計算にはPythonのLevenshteinライブラリを用いた。

実行結果画面の例

キュレーションされた和訳との差分の可視化

キュレーションされた和訳と各モデルの出力とのレーベンシュタイン距離



コマンドラインからAPIにアクセスし、遺伝子のSummaryを和訳する例

