

生成AIハンズオンセミナー(4)

Ollamaによるローカル環境での大規模言語モデルの実行

鳥井浩平

徳島大学デザイン型AI教育研究センター (併任)徳島大学理工学部



Ollama (オラマ) とは

- •大規模言語モデルをローカル環境で実行するためのソフトウェア
- Windows、macOS、Linux (Ubuntu) に対応



https://ollama.com/



ローカル環境とは





3



ローカル環境の利点・欠点

利点

- •実行時にインターネットの接続が不要
- •入力内容が外部に漏れる心配がない

<u>欠点</u>

- 高性能なPCでなければ、実行速度が遅い
- 最初は大規模言語モデルのダウンロードが必要
- 大規模言語モデルはサイズが大きく、PCの容量を圧迫する 1つのモデルで数ギガバイト



2024/11/12 生成AIハンズオンセミナー④

4

Ollamaをインストールする

※ここからはWindowsを前提として例示する



1 ウェブサイトにアクセス

https://ollama.com/



2024/11/12 生成AIハンズオンセミナー④

6

2 ダウンロード





3 インストーラを起動する

※ダブルクリックで実行

8





4 インストール①

육 Setup - Ollama version 0.3.14	- • ×
Ollama Windows Preview	23
Let's get you up and running with your own large language models.	}. <u>o</u> .{
Click Install to continue with the installation.	
	Install



5 インストール②

위 Setup - Ollama version 0.3.14	– 🗆 X
Installing Please wait while Setup installs Ollama on your computer.	([©])
Extracting files C:\Users\p1an0\AppData\Local\Programs\Ollama\lib\ollama\cublas64_11.dll	
	Cancel



6 コマンドプロンプトを開く



※「タスクバーにピン留め」しておくと後々便利



7 動作確認

C:¥User¥you> の右側に **ollama** と入力してEnterを押す
 下のような画面が出ればOK

C:¥Users¥p1a	nO> <mark>ollama</mark>
Usage:	ags]
ollama [fla	mmand]
Available Con serve create	nmanoj nmands: Start ollama Create a model from a Modelfile
show	Show information for a model
run	Run a model
stop	Stop a running model



Ollamaを使う

大規模言語モデル llama3.2 のダウンロードと実行



1 llama3.2 をダウンロード・実行する

- ① ollama run llama3.2 と入力してEnterを押す
- ② 下のような画面が出ればOK
- ※ 最初はモデルのダウンロードが入るため、時間がかかります





2 llama3.2 と会話する

※ここはChatGPT等と使い方は同じ

① >>> の右側に何か入力してEnterを押す

② しばらくすると llama3.2 から返答が返ってくる

>>> こんにちは! こんにちは!どういたしまして?何か相談したいことがありますか?

>>> 今日の晩御飯は何がいいかな? 今日の晩御飯を選ぶのはとても難しいことですね!

以下是一些人気の料理を紹介します。

ラーメン:たんぱくや牛肉、ひもの種子などを加えた味わいが良いでしょう。
 スープ:コップごとの熱いスープとハム、トマトなどを楽しむのも良いでしょう。
 ウドン: cold noodles で、タマゴ、ニンニク、パêsの種などを加えたみそ汁がอร่อยです。
 セリフライ:黄色くつったシダをコップごとのスープに染めたもので、甘味が強いです。

おすすめランキングは、これから教えますか?

>>> Send a message (/? for help)

3 会話を終了する

① /bye と入力(またはCtrlとDを同時に押す)と会話が終了する





いろんなモデルを使う

日本語特化モデルである llama3-elyza-jp-8b を使ってみる



1 モデルを検索する

Q Search models



Get up and running with large language models.

Run <u>Llama 3.2</u>, <u>Phi 3</u>, <u>Mistral</u>, <u>Gemma 2</u>, and other models. Customize and create your own.

Download ↓

Available for macOS, Linux, and Windows (preview)

今回は llama3 jp と入力して検索

Q Ilama3 jp

dsasai/llama3-elyza-jp-8b

stfate/llama3-elyza-jp-8b

pakachan/elyza-llama3-8b Llama-3-ELYZA-JP-8B is a large language model trained by ELYZA, Inc.

lucas2024/llama-3-elyza-jp-8b

jmorgan/llama3 A new model I just published!

View all \rightarrow

Run Llama 3.2. Phi 3. Mistral, Gemma 2, and

18



2 コマンドをコピーする

lsasai / I	lama3-elyza-jp-8b		
lo model sum	mary		
243 Pulls	Oupdated 3 months ago		
latest	 ✓ 1 Tag 	ollama run dsasai/llama3-elyza	-jp-8b
Updated 3 mont	hs ago	ecfdd9	2e89f6 · 4.9GB
model	arch llama · parameters 8.03B · quantization Q4_K_M		4.9GB
params	{ "stop": ["< start_header_id >", "< end_header_id >", "< eot_i…		128B
system	あなたは誠実で優秀な日本人のアシスタントです。特に指示が無い場合は、常		147B
	<pre>{{ if .System }}< start_header_id >system< end_header_id > {{ .S</pre>		



3 コマンドプロンプトに張り付けてEnter

※貼り付けはコマンドプロンプト上で右クリック





モデルの一覧表示と削除

PCの容量がいっぱいになる前に、いらないモデルは削除しよう



L ダウンロードしたモデルの一覧表示

- ① **ollama ls** と入力してEnterを押す
- ② ダウンロードしたモデルの一覧と情報が表示される
- ※ NAMEの列に書かれているのがモデルの名前





2 モデルを削除する

ollama rm モデルの名前 と入力してEnterを押す deleted モデルの名前 が表示されたら成功

C:¥Users¥p1an0>ollama ls NAME dsasai/llama3-elyza-jp-8b:latest llama3.2:latest

ID latest ecfdd92e89f6 a80c4f17acd5 SIZE

4.9 GB

2.0 GB

C:¥Users¥p1anO>ollama rm_llama3.2:latest deleted 'llama3.2:latest'



2024/11/12 生成AIハンズオンセミナー④

MODIFIED

6 minutes ago

46 minutes ago

さいごに



(参考) もう少しリッチな画面が好きな方は LM Studio



https://lmstudio.ai/

※使い方は下記ウェブサイトなどを参考に https://note.com/masayuki_abe/n/nd65ed694eec0

