第2章

## ChatGPT/Gemini/Claudeの基本コマンドから リアルタイム音声会話プログラム作成まで

# 独自 AI アプリ作りに役立つ LLM 用 API の使い方

村井 和夫

本稿では、ChatGPTに代表される大規模言語モデル (LLM: Large Language Model) を、プログラムに取り込んで使うための APIの利用法について次の2つを説明します。

- ①基本…curlコマンドとPythonコードからの利用 方法
- ②応用例…リアルタイム音声会話プログラム

APIを使ってプログラムに組み込むことによって、LLMを独自に構築しなくても、RAG (Retrieval Augmented Generation、検索拡張生成)やファインチューニングを使って特定目的に特化したLLMを作ったり、さまざまな独自の業務用アプリケーションを作ったりすることが可能になります。

### APIを利用するときに出てくる 基本用語

現在多数のLLMが公開されていますが、プログラムに取り込む技術は、従来からあるWeb APIと同じで、ネットワーク・プロトコルを通して使います、Web APIの使い方は非常に簡単ですぐ試すことができますが、基本的な用語を理解しておくとプログラムに組み込むときに役立ちます。

#### ● ネットワーク・プロトコルとエンド・ポイント

LLMのAPIも基本的に従来からあるWeb APIと同じです。Web APIで使われるネットワーク・プロトコルには単純なHTTP/HTTPSだけでなく、REST、SOAPなどいろいろなフレームワークがあります。現在では、多くの場合APIはHTTPSをベースにすることが多くなっています。しかし、リアルタイムの音声や映像を扱う必要がある場合、HTTPSを拡張したWebSocketや、リアルタイム通信専用のWebRTCのようなストリーミングに適したプロトコルを使用します。エンド・ポイントは、API機能を提供する接続先URLのことを指します。

#### ■ API キー

Web APIの基本は従量制課金なので、利用者の認

証が必要です.各社はAPIキーを発行して、それを使って利用者をネットワーク・プロトコルで認証して接続を許可します.そのため、ネットワーク・プロトコルは、必ずトランスポート・セキュリティに対応したプロトコルが使われます.このように、ネットワーク上でのセキュリティは担保されていますが、プログラムの中に直接書き込まないなど、利用する端末やプログラムで厳重に管理する必要があります.

#### ● コマンド・プロンプト

LLMに何を求めるかの質問です。プロンプト・エンジニアリングという用語があるように、求める役割と目的を明確に記述しないと、的確なレスポンスを得ることができません。基本的な役割や言語などの設定は、質問のコマンド・プロンプト以外に、別途設定するようになっている場合もあります。

#### ● コマンド・レスポンス形式

Web APIのコマンド・レスポンスの形式にはXMLやJSON形式などいろいろあります。現在ではほとんどがJSON形式になっています。従って、APIマニュアルを参考にして、JSON形式のコマンドを組み立てて送信し、JSON形式のレスポンスを受けて必要な情報を取り出す処理がプログラムの基本となります。

## 基本…curlコマンド or Pythonで 呼び出す(Gemini, ChatGPT, Claude)

代表的AIとして**表1**に示すモデルの使い方を紹介します.

いずれも、非常に安価なものからプログラムや理系に強い最新の高機能モデルまで実にさまざまなモデルがありますが、執筆時点で無料枠が一番大きく安価なのはGeminiでした。今回は、AIのAPIの利用方法を確認することを目的として、執筆時点で一番安価で簡単な会話レベルで利用しやすいモデルを選び、APIの利用方法を紹介します。