# ■ Internship Application Manager
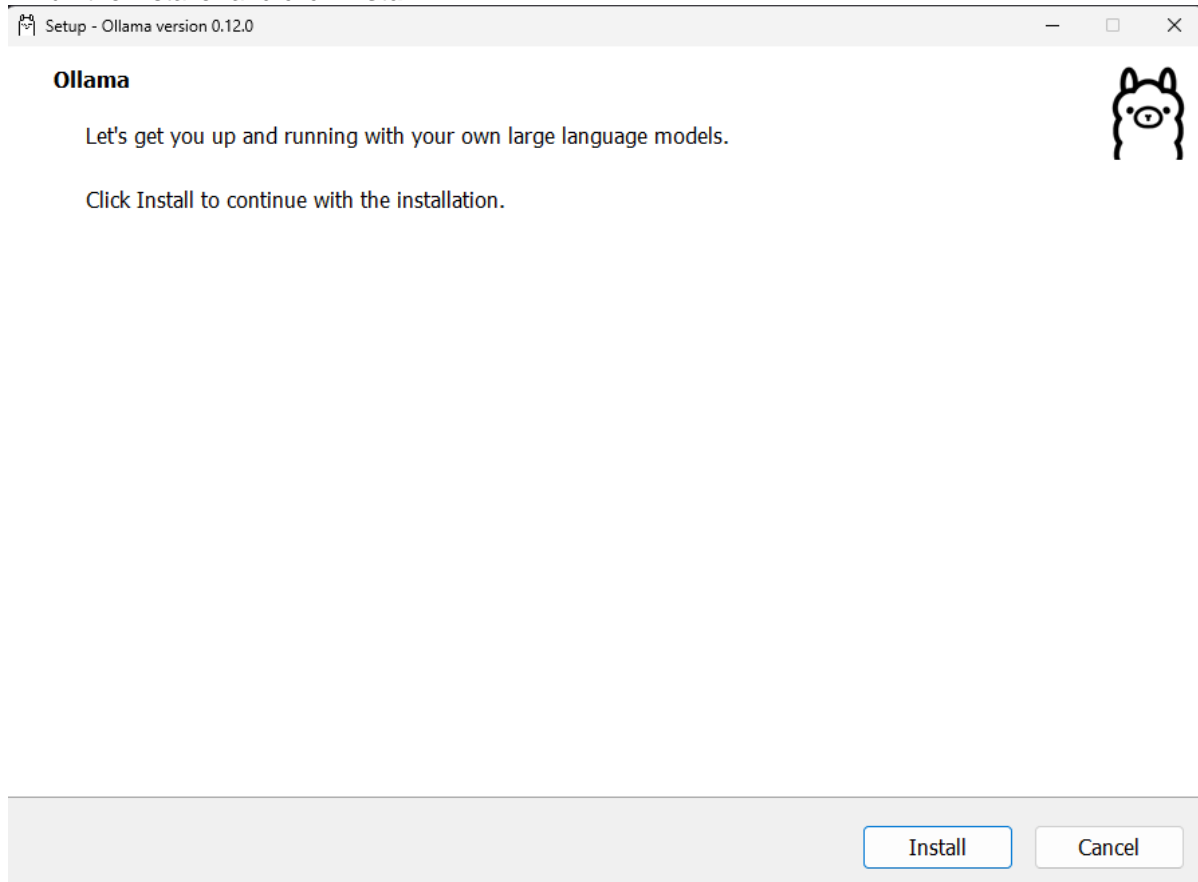
## Ollama AI Setup Guide (Windows)

This guide will walk you through installing and setting up Ollama so you can use the AI Assistant in the Internship Application Manager. Follow each step carefully and you'll be ready to go in minutes.

Prepared for: Project Team

Version 1.0

## Step 1: Install Ollama

1. Download the Ollama installer for Windows from https://ollama.com/download

2. Run the installer and click **Install**:



## Step 2: Start Ollama

You can start Ollama in two ways:

• **Option 1:** Open the Ollama app from the Start Menu.
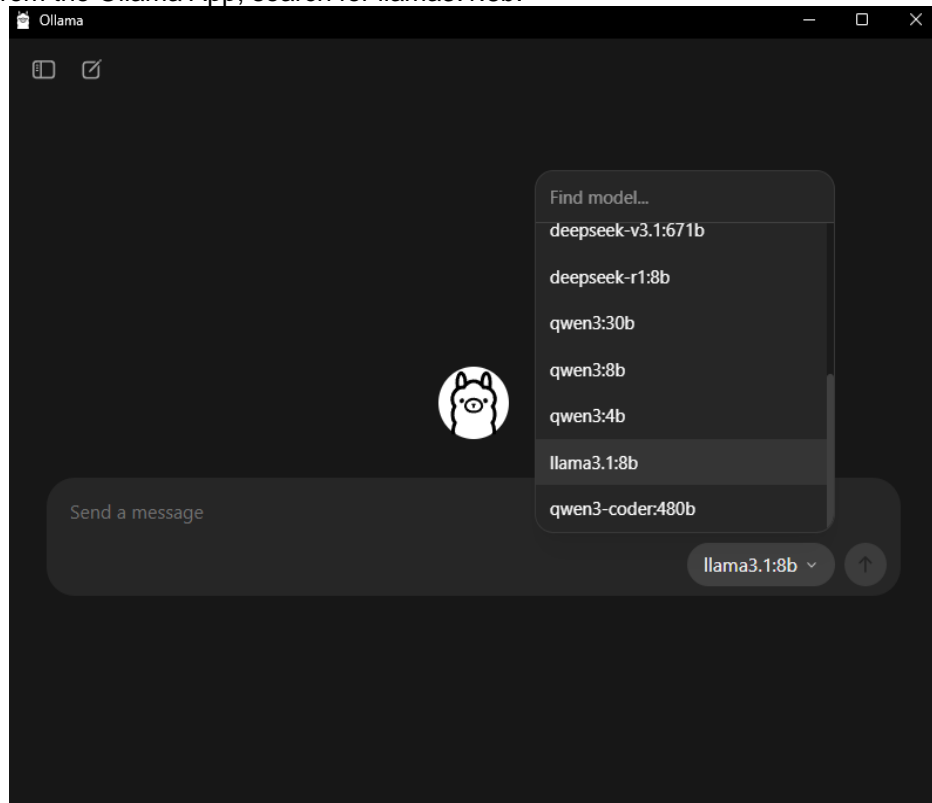
• **Option 2:** Use Command Prompt and run:

```
ollama serve
```

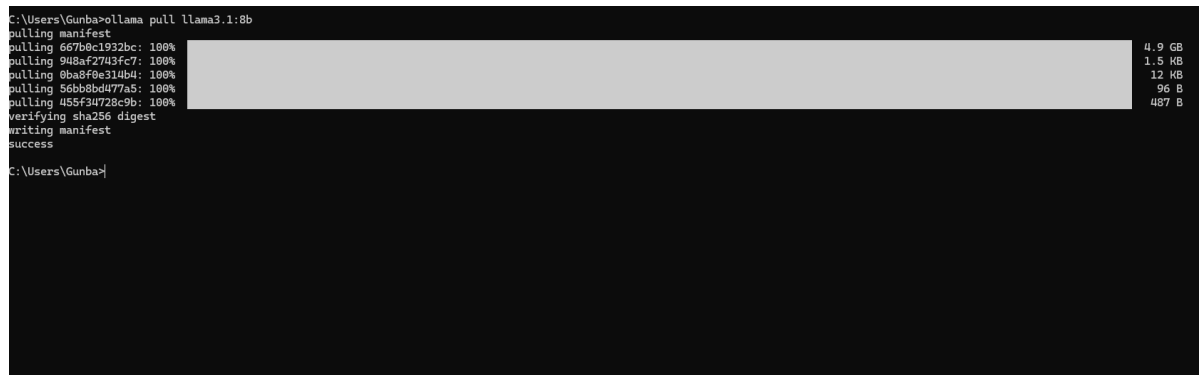## Step 3: Download the Required Model

The Internship Application Manager uses the **llama3.1:8b** model.

Option 1: From the Ollama App, search for llama3.1:8b:



Option 2: From Command Prompt, run:

```
ollama pull llama3.1:8b
```



## Running the Internship Application Manager

■■ Before running the app, make sure Ollama is running — either the app is open or you see it in Task Manager (Windows).

From the repository root, run:

```
dotnet run
```

This will start the web app locally. Open http://localhost:5000/ in your browser.

## Troubleshooting

• If you see 'server not running' or 'connection refused', make sure Ollama is running and the model is downloaded.

• Restart both Ollama and your backend if needed.

• Check firewall/antivirus settings.

• Run 'ollama serve' manually to view error messages.

## ■ Quick Checklist

| | |
|---|---|
| ✔ | Ollama is installed |
| ✔ | Ollama is running (app or ollama serve) |
| ✔ | Model llama3.1:8b is downloaded |
| ✔ | Backend is running (dotnet run) |
| ✔ | Access http://localhost:5000/ in your browser |

## ■ All Done!

If you followed all the steps, hooray! Your Internship Application Manager with Ollama AI is ready to go. Enjoy tailoring resumes and exploring internship postings!