

# What is OpenRAG?

OpenRAG is an open-source package for building agentic RAG systems. It supports integration with a wide range of orchestration tools, vector databases, and LLM providers.

OpenRAG connects and amplifies three popular, proven open-source projects into one powerful platform:

- [Langflow](#) - Langflow is a powerful tool to build and deploy AI agents and MCP servers. It supports all major LLMs, vector databases and a growing library of AI tools.
- [OpenSearch](#) - OpenSearch is a community-driven, Apache 2.0-licensed open source search and analytics suite that makes it easy to ingest, search, visualize, and analyze data.
- [Docling](#) - Docling simplifies document processing, parsing diverse formats — including advanced PDF understanding — and providing seamless integrations with the gen AI ecosystem.

OpenRAG builds on Langflow's familiar interface while adding OpenSearch for vector storage and Docling for simplified document parsing, with opinionated flows that serve as ready-to-use recipes for ingestion, retrieval, and generation from popular sources like OneDrive, Google Drive, and AWS.

What's more, every part of the stack is swappable. Write your own custom components in Python, try different language models, and customize your flows to build an agentic RAG system.

Ready to get started? [Install OpenRAG](#) and then run the [Quickstart](#) to create a powerful RAG pipeline.

# Install OpenRAG

Install the [OpenRAG Python wheel](#), and then run the [OpenRAG Terminal User Interface\(TUI\)](#) to start your OpenRAG deployment with a guided setup process.

If you prefer running Docker commands and manually editing `.env` files, see [Deploy with Docker](#).

## Prerequisites

- [Python Version 3.10 to 3.13](#)
- [uv](#)
- [Podman](#) (recommended) or [Docker](#) installed
- [Docker Compose](#) installed. If using Podman, use `podman-compose` or alias Docker compose commands to Podman commands.
- Create an [OpenAI API key](#). This key is **required** to start OpenRAG, but you can choose a different model provider during [Application Onboarding](#).
- Optional: GPU support requires an NVIDIA GPU with [CUDA](#) support and compatible NVIDIA drivers installed on the OpenRAG host machine. If you don't have GPU capabilities, OpenRAG provides an alternate CPU-only deployment.

## Install the OpenRAG Python wheel

### ⚠ IMPORTANT

The `.whl` file is currently available as an internal download during public preview, and will be published to PyPI in a future release.

The OpenRAG wheel installs the Terminal User Interface (TUI) for configuring and running OpenRAG.

1. Create a new project with a virtual environment using `uv init`.

```
uv init YOUR_PROJECT_NAME
cd YOUR_PROJECT_NAME
```

The `(venv)` prompt doesn't change, but `uv` commands will automatically use the project's virtual environment. For more information on virtual environments, see the

[uv documentation](#).

2. Add the local OpenRAG wheel to your project's virtual environment.

```
uv add PATH/T0/openrag-VERSION-py3-none-any.whl
```

Replace `PATH/T0/` and `VERSION` with the path and version of your downloaded OpenRAG `.whl` file.

For example, if your `.whl` file is in the `~/Downloads` directory, the command is `uv add ~/Downloads/openrag-0.1.8-py3-none-any.whl`.

3. Ensure all dependencies are installed and updated in your virtual environment.

```
uv sync
```

4. Start the OpenRAG TUI.

```
uv run openrag
```

5. Continue with [Setup OpenRAG with the TUI](#).

## Set up OpenRAG with the TUI

The TUI creates a `.env` file in your OpenRAG directory root and starts OpenRAG. If the TUI detects a `.env` file in the OpenRAG root directory, it sources any variables from the `.env` file. If the TUI detects OAuth credentials, it enforces the **Advanced Setup** path.

**Basic Setup** generates all of the required values for OpenRAG except the OpenAI API key. **Basic Setup** does not set up OAuth connections for ingestion from cloud providers. For OAuth setup, use **Advanced Setup**.

**Basic Setup** and **Advanced Setup** enforce the same authentication settings for the Langflow server, but manage document access differently. For more information, see [Authentication and document access](#).

## Basic setup

## Advanced setup

1. To install OpenRAG with **Basic Setup**, click **Basic Setup** or press `1`.
2. Click **Generate Passwords** to generate passwords for OpenSearch and Langflow.
3. Paste your OpenAI API key in the OpenAI API key field.
4. Click **Save Configuration**.
5. To start OpenRAG, click **Start Container Services**. Startup pulls container images and runs them, so it can take some time. When startup is complete, the TUI displays the following:

```
Services started successfully  
Command completed successfully
```

6. To open the OpenRAG application, click **Open App**.
7. Continue with [Application Onboarding](#).

## Application onboarding

The first time you start OpenRAG, whether using the TUI or a `.env` file, you must complete application onboarding.

Most values from onboarding can be changed later in the OpenRAG **Settings** page, but there are important restrictions.

The **language model provider** and **embeddings model provider** can only be selected at onboarding, and you must use the same provider for your language model and embedding model. To change your provider selection later, you must completely reinstall OpenRAG.

The **language model** can be changed later in **Settings**, but the **embeddings model** cannot be changed later.

## OpenAI

## IBM watsonx.ai

## Ollama

1. Enable **Get API key from environment variable** to automatically enter your key from the TUI-generated `.env` file.
2. Under **Advanced settings**, select your **Embedding Model** and **Language Model**.
3. To load 2 sample PDFs, enable **Sample dataset**. This is recommended, but not required.
4. Click **Complete**.
5. Continue with the [Quickstart](#).

# Deploy with Docker

There are two different Docker Compose files. They deploy the same applications and containers, but to different environments.

- `docker-compose.yml` is an OpenRAG deployment with GPU support for accelerated AI processing.
- `docker-compose-cpu.yml` is a CPU-only version of OpenRAG for systems without GPU support. Use this Docker compose file for environments where GPU drivers aren't available.

Both Docker deployments depend on `docling serve` to be running on port `5001` on the host machine. This enables [Mac MLX](#) support for document processing. Installing OpenRAG with the TUI starts `docling serve` automatically, but for a Docker deployment you must manually start the `docling serve` process.

## Prerequisites

- [Python Version 3.10 to 3.13](#)
- [uv](#)
- [Podman](#) (recommended) or [Docker](#) installed
- [Docker Compose](#) installed. If you're using Podman, use `podman-compose` or alias Docker compose commands to Podman commands.
- Create an [OpenAI API key](#). This key is **required** to start OpenRAG, but you can choose a different model provider during [Application Onboarding](#).
- Optional: GPU support requires an NVIDIA GPU with CUDA support and compatible NVIDIA drivers installed on the OpenRAG host machine. If you don't have GPU capabilities, OpenRAG provides an alternate CPU-only deployment.

## Deploy OpenRAG with Docker Compose

To install OpenRAG with Docker Compose, do the following:

1. Clone the OpenRAG repository.

```
git clone https://github.com/langflow-ai/openrag.git
cd openrag
```

2. Install dependencies.

```
uv sync
```

3. Copy the example `.env` file included in the repository root. The example file includes all environment variables with comments to guide you in finding and setting their values.

```
cp .env.example .env
```

Alternatively, create a new `.env` file in the repository root.

```
touch .env
```

4. Set environment variables. The Docker Compose files will be populated with values from your `.env`. The following values are **required** to be set:

```
OPENSEARCH_PASSWORD=your_secure_password
OPENAI_API_KEY=your_openai_api_key
LANGFLOW_SUPERUSER=admin
LANGFLOW_SUPERUSER_PASSWORD=your_langflow_password
LANGFLOW_SECRET_KEY=your_secret_key
```

For more information on configuring OpenRAG with environment variables, see [Environment variables](#).

5. Start `docling serve` on the host machine. Both Docker deployments depend on `docling serve` to be running on port `5001` on the host machine. This enables [Mac MLX](#) support for document processing.

```
uv run python scripts/docling_ctl.py start --port 5001
```

6. Confirm `docling serve` is running.

```
uv run python scripts/docling_ctl.py status
```

Successful result:

```
Status: running
Endpoint: http://127.0.0.1:5001
Docs: http://127.0.0.1:5001/docs
PID: 27746
```

7. Deploy OpenRAG with Docker Compose based on your deployment type.

For GPU-enabled systems, run the following commands:

```
docker compose build
docker compose up -d
```

For environments without GPU support, run:

```
docker compose -f docker-compose-cpu.yml up -d
```

The OpenRAG Docker Compose file starts five containers:

Container Name	Default Address	Purpose
OpenRAG Backend	<a href="http://localhost:8000">http://localhost:8000</a>	FastAPI server and core functionality.
OpenRAG Frontend	<a href="http://localhost:3000">http://localhost:3000</a>	React web interface for users.
Langflow	<a href="http://localhost:7860">http://localhost:7860</a>	AI workflow engine and flow management.
OpenSearch	<a href="http://localhost:9200">http://localhost:9200</a>	Vector database for document storage.
OpenSearch Dashboards	<a href="http://localhost:5601">http://localhost:5601</a>	Database administration interface.

8. Verify installation by confirming all services are running.



```
docker compose ps
```

You can now access the application at:

- **Frontend:** <http://localhost:3000>
- **Backend API:** <http://localhost:8000>
- **Langflow:** <http://localhost:7860>

9. Continue with [Application Onboarding](#).

To stop `docling serve` when you're done with your OpenRAG deployment, run:

```
uv run python scripts/docling_ctl.py stop
```

## Application onboarding

The first time you start OpenRAG, whether using the TUI or a `.env` file, you must complete application onboarding.

Most values from onboarding can be changed later in the OpenRAG **Settings** page, but there are important restrictions.

The **language model provider** and **embeddings model provider** can only be selected at onboarding, and you must use the same provider for your language model and embedding model. To change your provider selection later, you must completely reinstall OpenRAG.

The **language model** can be changed later in **Settings**, but the **embeddings model** cannot be changed later.

**OpenAI**    **IBM watsonx.ai**    **Ollama**

1. Enable **Get API key from environment variable** to automatically enter your key from the TUI-generated `.env` file.
2. Under **Advanced settings**, select your **Embedding Model** and **Language Model**.

3. To load 2 sample PDFs, enable **Sample dataset**. This is recommended, but not required.
4. Click **Complete**.
5. Continue with the [Quickstart](#).

## Container management commands

Manage your OpenRAG containers with the following commands. These commands are also available in the TUI's [Status menu](#).

### Upgrade containers

Upgrade your containers to the latest version while preserving your data.

```
docker compose pull
docker compose up -d --force-recreate
```

### Rebuild containers (destructive)

Reset state by rebuilding all of your containers. Your OpenSearch and Langflow databases will be lost. Documents stored in the `./documents` directory will persist, since the directory is mounted as a volume in the OpenRAG backend container.

```
docker compose up --build --force-recreate --remove-orphans
```

### Remove all containers and data (destructive)

Completely remove your OpenRAG installation and delete all data. This deletes all of your data, including OpenSearch data, uploaded documents, and authentication.

```
docker compose down --volumes --remove-orphans --rmi local
docker system prune -f
```



# Quickstart

Get started with OpenRAG by loading your knowledge, swapping out your language model, and then chatting with the OpenRAG API.


## Prerequisites

- [Install and start OpenRAG](#)

## Find your way around


1. In OpenRAG, click  **Chat**. The chat is powered by the OpenRAG OpenSearch Agent. For more information, see [Langflow Agents](#).
2. Ask `What documents are available to you?` The agent responds with a message summarizing the documents that OpenRAG loads by default, which are PDFs about evaluating data quality when using LLMs in health care. Knowledge is stored in OpenSearch. For more information, see [Knowledge](#).
3. To confirm the agent is correct, click  **Knowledge**. The **Knowledge** page lists the documents OpenRAG has ingested into the OpenSearch vector database. Click on a document to display the chunks derived from splitting the default documents into the vector database.

## Add your own knowledge

1. To add documents to your knowledge base, click  **Add Knowledge**.
  - Select **Add File** to add a single file from your local machine (mapped with the Docker volume mount).
  - Select **Process Folder** to process an entire folder of documents from your local machine (mapped with the Docker volume mount).
  - Select your cloud storage provider to add knowledge from an OAuth-connected storage provider. For more information, see [OAuth ingestion](#).
2. Return to the Chat window and ask a question about your loaded data. For example, with a manual about a PC tablet loaded, ask `How do I connect this device to WiFi?` The agent responds with a message indicating it now has your knowledge as context for answering questions.
3. Click the **Function Call: search\_documents (tool\_call)** that is printed in the Playground. These events log the agent's request to the tool and the tool's

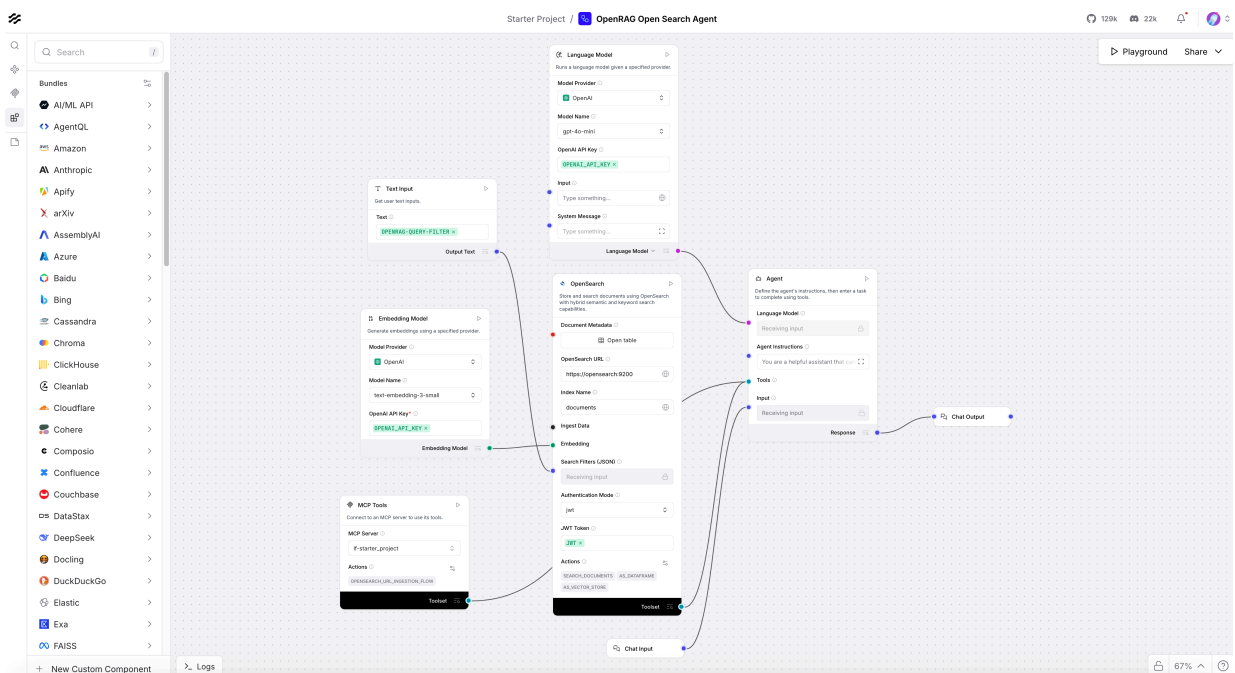
response, so you have direct visibility into your agent's functionality. If you aren't getting the results you need, you can further tune the knowledge ingestion and agent behavior in the next section.

## Swap out the language model to modify agent behavior

To modify the knowledge ingestion or Agent behavior, click  **Settings**.

In this example, you'll try a different LLM to demonstrate how the Agent's response changes. You can only change the **Language model**, and not the **Model provider** that you started with in OpenRAG. If you're using Ollama, you can use any installed model.

1. To edit the Agent's behavior, click **Edit in Langflow**. You can more quickly access the **Language Model** and **Agent Instructions** fields in this page, but for illustration purposes, navigate to the Langflow visual builder.
2. OpenRAG warns you that you're entering Langflow. Click **Proceed**.
3. The OpenRAG OpenSearch Agent flow appears.



4. In the **Language Model** component, under **Model**, select a different OpenAI model.
5. Save your flow with `Command+S`.
6. In OpenRAG, start a new conversation by clicking the **+** in the **Conversations** tab.


7. Ask the same question as before to demonstrate how a different language model changes the results.

## Integrate OpenRAG into your application

To integrate OpenRAG into your application, use the [Langflow API](#). Make requests with Python, TypeScript, or any HTTP client to run one of OpenRAG's default flows and get a response, and then modify the flow further to improve results. Langflow provides code snippets to help you get started.

1. Create a [Langflow API key](#).

▶ [Create a Langflow API key](#)

2. To navigate to the OpenRAG OpenSearch Agent flow, click  **Settings**, and then click **Edit in Langflow** in the OpenRAG OpenSearch Agent flow.
3. Click **Share**, and then click **API access**.

The default code in the API access pane constructs a request with the Langflow server `url`, `headers`, and a `payload` of request data. The code snippets automatically include the `LANGFLOW_SERVER_ADDRESS` and `FLOW_ID` values for the flow. Replace these values if you're using the code for a different server or flow. The default Langflow server address is <http://localhost:7860>.

**Python**    TypeScript    curl

```
import requests
import os
import uuid
api_key = 'LANGFLOW_API_KEY'
url = "http://LANGFLOW_SERVER_ADDRESS/api/v1/run/FLOW_ID"
# The complete API endpoint URL for this flow
# Request payload configuration
payload = {
    "output_type": "chat",
    "input_type": "chat",
```

```

    "input_value": "hello world!"
}
payload["session_id"] = str(uuid.uuid4())
headers = {"x-api-key": api_key}
try:
    # Send API request
    response = requests.request("POST", url, json=payload,
headers=headers)
    response.raise_for_status() # Raise exception for bad
status codes
    # Print response
    print(response.text)
except requests.exceptions.RequestException as e:
    print(f"Error making API request: {e}")
except ValueError as e:
    print(f"Error parsing response: {e}")

```

4. Copy the snippet, paste it in a script file, and then run the script to send the request. If you are using the curl snippet, you can run the command directly in your terminal.

If the request is successful, the response includes many details about the flow run, including the session ID, inputs, outputs, components, durations, and more. The following is an example of a response from running the **Simple Agent** template flow:

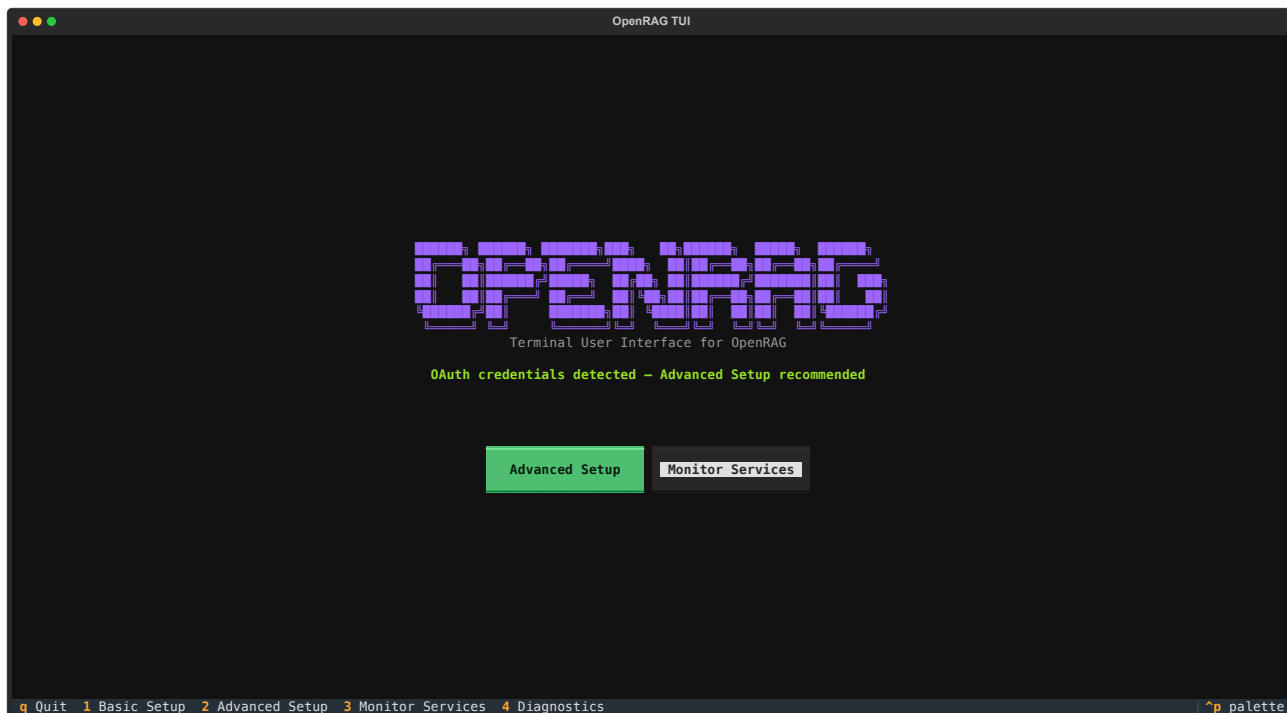
▶ Result

To further explore the API, see:

- The Langflow [Quickstart](#) extends this example with extracting fields from the response.
- [Get started with the Langflow API](#)

# Terminal User Interface (TUI) commands

The OpenRAG Terminal User Interface (TUI) allows you to set up, configure, and monitor your OpenRAG deployment directly from the terminal, on any operating system.



Instead of starting OpenRAG using Docker commands and manually editing values in the `.env` file, the TUI walks you through the setup. It prompts for variables where required, creates a `.env` file for you, and then starts OpenRAG.

Once OpenRAG is running, use the TUI to monitor your application, control your containers, and retrieve logs.

## Start the TUI

To start the TUI, run the following commands from the directory where you installed OpenRAG.

```
uv sync
uv run openrag
```

The TUI Welcome Screen offers basic and advanced setup options. For more information on setup values during installation, see [Install OpenRAG](#).

# Navigation

The TUI accepts mouse input or keyboard commands.

- `Arrow keys`: move between options
- `Tab / Shift+Tab`: switch fields and buttons
- `Enter`: select/confirm
- `Escape`: back
- `Q`: quit
- `Number keys (1-4)`: quick access to main screens

## Container management

The TUI can deploy, manage, and upgrade your OpenRAG containers.

### Start container services

Click **Start Container Services** to start the OpenRAG containers. The TUI automatically detects your container runtime, and then checks if your machine has compatible GPU support by checking for `CUDA`, `NVIDIA_SMI`, and Docker/Podman runtime support. This check determines which Docker Compose file OpenRAG uses. The TUI then pulls the images and deploys the containers with the following command.

```
docker compose up -d
```

If images are missing, the TUI runs `docker compose pull`, then runs `docker compose up -d`.

### Start native services

A "native" service in OpenRAG refers to a service run natively on your machine, and not within a container. The `docling serve` process is a native service in OpenRAG, because it's a document processing service that is run on your local machine, and controlled separately from the containers.

To start or stop `docling serve` or any other native services, in the TUI main menu, click **Start Native Services** or **Stop Native Services**.

To view the status, port, or PID of a native service, in the TUI main menu, click [Status](#).



## Status

The **Status** menu displays information on your container deployment. Here you can check container health, find your service ports, view logs, and upgrade your containers.

To view streaming logs, select the container you want to view, and press `l`. To copy your logs, click **Copy to Clipboard**.

To **upgrade** your containers, click **Upgrade**. **Upgrade** runs `docker compose pull` and then `docker compose up -d --force-recreate`. The first command pulls the latest images of OpenRAG. The second command recreates the containers with your data persisted.

To **reset** your containers, click **Reset**. Reset gives you a completely fresh start. Reset deletes all of your data, including OpenSearch data, uploaded documents, and authentication. **Reset** runs two commands. It first stops and removes all containers, volumes, and local images.

```
docker compose down --volumes --remove-orphans --rmi local
```

When the first command is complete, OpenRAG removes any additional Docker objects with `prune`.

```
docker system prune -f
```

## Diagnostics

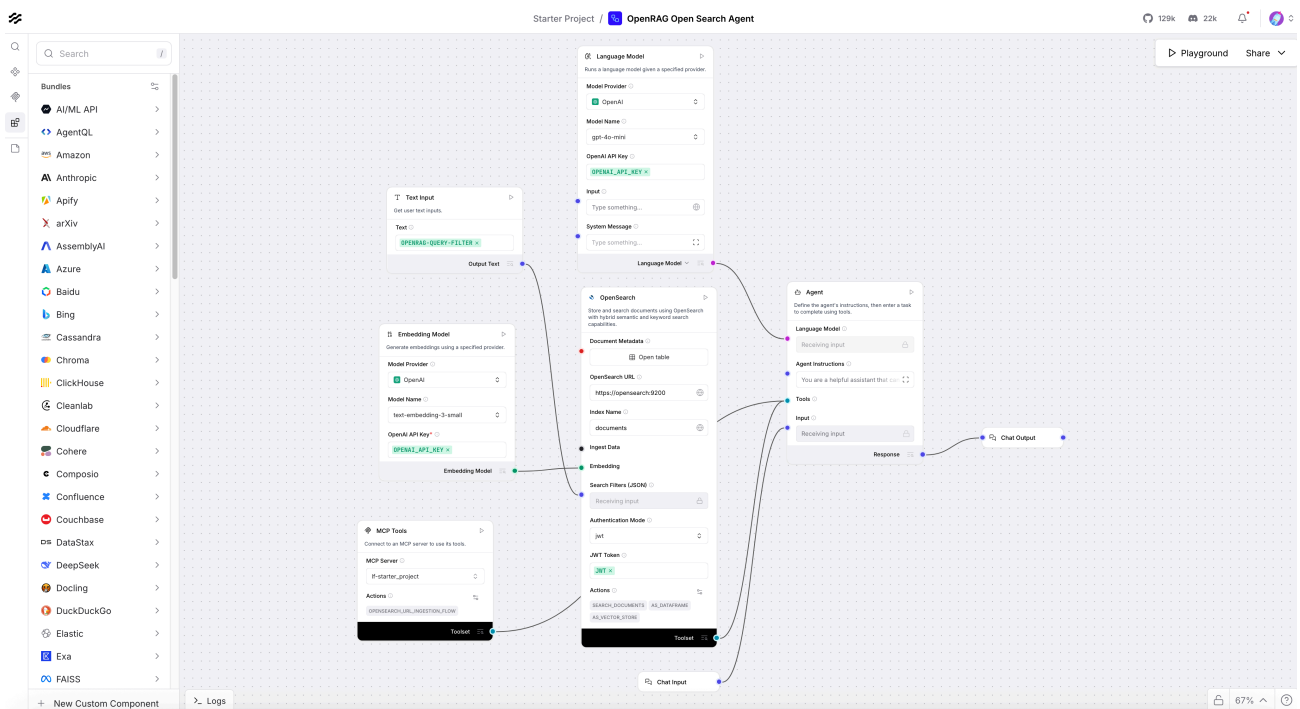
The **Diagnostics** menu provides health monitoring for your container runtimes and monitoring of your OpenSearch security.

# Langflow Agents

OpenRAG leverages Langflow's Agent component to power the OpenRAG OpenSearch Agent flow.

Flows in Langflow are functional representations of application workflows, with multiple component nodes connected as single steps in a workflow.

In the OpenRAG OpenSearch Agent flow, components like the Langflow **Agent component** and **OpenSearch component** are connected to intelligently chat with your knowledge by embedding your query, comparing it the vector database embeddings, and generating a response with the LLM.




The Agent component shines here in its ability to make decisions on not only what query should be sent, but when a query is necessary to solve the problem at hand.


▶ How do agents work?

## Use the OpenRAG OpenSearch Agent flow


If you've chatted with your knowledge in OpenRAG, you've already experienced the OpenRAG OpenSearch Agent chat flow. To switch OpenRAG over to the **Langflow visual**

[editor](#) and view the OpenRAG OpenSearch Agentflow, click  **Settings**, and then click **Edit in Langflow**. This flow contains eight components connected together to chat with your data:

- The **Agent component** orchestrates the entire flow by deciding when to search the knowledge base, how to formulate search queries, and how to combine retrieved information with the user's question to generate a comprehensive response. The **Agent** behaves according to the prompt in the **Agent Instructions** field.
- The **Chat Input component** is connected to the Agent component's Input port. This allows the flow to be triggered by an incoming prompt from a user or application.
- The **OpenSearch component** is connected to the Agent component's Tools port. The agent may not use this database for every request; the agent only uses this connection if it decides the knowledge can help respond to the prompt.
- The **Language Model component** is connected to the Agent component's Language Model port. The agent uses the connected LLM to reason through the request sent through Chat Input.
- The **Embedding Model component** is connected to the OpenSearch component's Embedding port. This component converts text queries into vector representations that are compared with document embeddings stored in OpenSearch for semantic similarity matching. This gives your Agent's queries context.
- The **Text Input component** is populated with the global variable `OPENRAG-QUERY-FILTER`. This filter is the **Knowledge filter**, and filters which knowledge sources to search through.
- The **Agent component's** Output port is connected to the **Chat Output component**, which returns the final response to the user or application.
- An **MCP Tools component** is connected to the Agent's **Tools** port. This component calls the **OpenSearch URL Ingestion flow**, which Langflow uses as an MCP server to fetch content from URLs and store in OpenSearch.

All flows included with OpenRAG are designed to be modular, performant, and provider-agnostic. To modify a flow, click  **Settings**, and click **Edit in Langflow**. OpenRAG's visual editor is based on the [Langflow visual editor](#), so you can edit your flows to match your specific use case.

For an example of changing out the agent's language model in OpenRAG, see the [Quickstart](#).

To restore the flow to its initial state, in OpenRAG, click  **Settings**, and then click **Restore Flow**. OpenRAG warns you that this discards all custom settings. Click **Restore** to restore the flow.

## Additional Langflow functionality

Langflow includes features beyond Agents to help you integrate OpenRAG into your application, and all Langflow features are included in OpenRAG.

- Langflow can serve your flows as an [MCP server](#), or consume other MCP servers as an [MCP client](#). Get started with the [MCP tutorial](#).
- If you don't see the component you need, extend Langflow's functionality by creating [custom Python components](#).
- Langflow offers component [bundles](#) to integrate with many popular vector stores, AI/ML providers, and search APIs.

# OpenSearch Knowledge

OpenRAG uses [OpenSearch](#) for its vector-backed knowledge store. This is a specialized database for storing and retrieving embeddings, which helps your Agent efficiently find relevant information. OpenSearch provides powerful hybrid search capabilities with enterprise-grade security and multi-tenancy support.

## Authentication and document access

OpenRAG supports two authentication modes based on how you [install OpenRAG](#), and which mode you choose affects document access.

**No-auth mode (Basic Setup):** This mode uses a single anonymous JWT token for OpenSearch authentication, so documents uploaded to the `documents` index by one user are visible to all other users on the OpenRAG server.

**OAuth mode (Advanced Setup):** Each OpenRAG user is granted a JWT token, and each document is tagged with user ownership. Documents are filtered by user ownership, ensuring users only see documents they uploaded or have access to.

## Ingest knowledge

OpenRAG supports knowledge ingestion through direct file uploads and OAuth connectors. To configure the knowledge ingestion pipeline parameters, see [Docling Ingestion](#).

### Direct file ingestion

The **Knowledge Ingest** flow uses Langflow's [File component](#) to split and embed files loaded from your local machine into the OpenSearch database.

The default path to your local folder is mounted from the `./documents` folder in your OpenRAG project directory to the `/app/documents/` directory inside the Docker container. Files added to the host or the container will be visible in both locations. To configure this location, modify the **Documents Paths** variable in either the TUI's [Advanced Setup](#) menu or in the `.env` used by Docker Compose.

To load and process a single file from the mapped location, click **+ Add Knowledge**, and then click **Add File**. The file is loaded into your OpenSearch database, and appears

in the Knowledge page.

To load and process a directory from the mapped location, click **+** **Add Knowledge**, and then click **Process Folder**. The files are loaded into your OpenSearch database, and appear in the Knowledge page.

## Ingest files through OAuth connectors

OpenRAG supports Google Drive, OneDrive, and AWS S3 as OAuth connectors for seamless document synchronization.

OAuth integration allows individual users to connect their personal cloud storage accounts to OpenRAG. Each user must separately authorize OpenRAG to access their own cloud storage files. When a user connects a cloud service, they are redirected to authenticate with that service provider and grant OpenRAG permission to sync documents from their personal cloud storage.

Before users can connect their cloud storage accounts, you must configure OAuth credentials in OpenRAG. This requires registering OpenRAG as an OAuth application with a cloud provider and obtaining client ID and secret keys for each service you want to support.

To add an OAuth connector to OpenRAG, do the following. This example uses Google OAuth. If you wish to use another provider, add the secrets to another provider.

### **TUI** `.env`

1. If OpenRAG is running, stop it with **Status > Stop Services**.
2. Click **Advanced Setup**.
3. Add the OAuth provider's client and secret key in the **Advanced Setup** menu.
4. Click **Save Configuration**. The TUI generates a new `.env` file with your OAuth values.
5. Click **Start Container Services**.

The OpenRAG frontend at `http://localhost:3000` now redirects to an OAuth callback login page for your OAuth provider. A successful authentication opens OpenRAG with the required scopes for your connected storage.

To add knowledge from an OAuth-connected storage provider, do the following:

1. Click **+** **Add Knowledge**, and then select the storage provider, for example, **Google Drive**. The **Add Cloud Knowledge** page opens.
2. To add files or folders from the connected storage, click **+** **Add Files**. Select the files or folders you want and click **Select**. You can select multiples.
3. When your files are selected, click **Ingest Files**. The ingestion process may take some time, depending on the size of your documents.
4. When ingestion is complete, your documents are available in the Knowledge screen.

## Explore knowledge

The **Knowledge** page lists the documents OpenRAG has ingested into the OpenSearch vector database's `documents` index.

To explore your current knowledge, click **||| Knowledge**. Click on a document to display the chunks derived from splitting the default documents into the vector database.

Documents are processed with the default **Knowledge Ingest** flow, so if you want to split your documents differently, edit the **Knowledge Ingest** flow.





All flows included with OpenRAG are designed to be modular, performant, and provider-agnostic. To modify a flow, click **⚙ Settings**, and click **Edit in Langflow**. OpenRAG's visual editor is based on the [Langflow visual editor](#), so you can edit your flows to match your specific use case.

## Create knowledge filters



OpenRAG includes a knowledge filter system for organizing and managing document collections. Knowledge filters are saved search configurations that allow you to create custom views of your document collection. They store search queries, filter criteria, and display settings that can be reused across different parts of OpenRAG.


Knowledge filters help agents work more efficiently with large document collections by focusing their context within relevant documents sets.

To create a knowledge filter, do the following:

1. Click  **All Knowledge**, and then click  **Create New Filter**. The **Create New Knowledge Filter** pane appears.
2. Enter a **Name** and **Description**, and then click  **Create Filter**. A new filter is created with default settings that match everything.
3. To modify the default filter, click  **All Knowledge**, and then click your new filter to edit it in the **Knowledge Filter** pane.

The following filter options are configurable.

- **Search Query:** Enter text for semantic search, such as "financial reports from Q4".
  - **Data Sources:** Select specific data sources or folders to include.
  - **Document Types:** Filter by file type.
  - **Owners:** Filter by who uploaded the documents.
  - **Sources:** Filter by connector types, such as local upload or Google Drive.
  - **Result Limit:** Set maximum number of results. The default is .
  - **Score Threshold:** Set minimum relevance score. The default score is .
4. When you're done editing the filter, click  **Save Configuration**.
  5. To apply the filter to OpenRAG globally, click  **All Knowledge**, and then select the filter to apply.

To apply the filter to a single chat session, in the  **Chat** window, click @, and then select the filter to apply.

## OpenRAG default configuration

OpenRAG automatically detects and configures the correct vector dimensions for embedding models, ensuring optimal search performance and compatibility.

The complete list of supported models is available at [models\\_service.py](#) in the [OpenRAG repository](#).

You can use custom embedding models by specifying them in your configuration.



If you use an unknown embedding model, OpenRAG will automatically fall back to 1536 dimensions and log a warning. The system will continue to work, but search quality may be affected if the actual model dimensions differ from 1536.

The default embedding dimension is 1536 and the default model is text-embedding-3-small.

For models with known vector dimensions, see [settings.py](#) in the OpenRAG repository.

# Docling Ingestion

OpenRAG uses **Docling** for its document ingestion pipeline. More specifically, OpenRAG uses **Docling Serve**, which starts a `docling serve` process on your local machine and runs Docling ingestion through an API service.

Docling ingests documents from your local machine or OAuth connectors, splits them into chunks, and stores them as separate, structured documents in the OpenSearch `documents` index.

OpenRAG chose Docling for its support for a wide variety of file formats, high performance, and advanced understanding of tables and images.

## Docling ingestion settings

These settings configure the Docling ingestion parameters.

OpenRAG will warn you if `docling serve` is not running. To start or stop `docling serve` or any other native services, in the TUI main menu, click **Start Native Services** or **Stop Native Services**.

**Embedding model** determines which AI model is used to create vector embeddings. The default is `text-embedding-3-small`.

**Chunk size** determines how large each text chunk is in number of characters. Larger chunks yield more context per chunk, but may include irrelevant information. Smaller chunks yield more precise semantic search, but may lack context. The default value of `1000` characters provides a good starting point that balances these considerations.

**Chunk overlap** controls the number of characters that overlap over chunk boundaries. Use larger overlap values for documents where context is most important, and use smaller overlap values for simpler documents, or when optimization is most important. The default value of 200 characters of overlap with a chunk size of 1000 (20% overlap) is suitable for general use cases. Decrease the overlap to 10% for a more efficient pipeline, or increase to 40% for more complex documents.

**OCR** enables or disabled OCR processing when extracting text from images and scanned documents. OCR is disabled by default. This setting is best suited for processing text-

based documents as quickly as possible with Docling's `DocumentConverter`. Images are ignored and not processed.

Enable OCR when you are processing documents containing images with text that requires extraction, or for scanned documents. Enabling OCR can slow ingestion performance.

If OpenRAG detects that the local machine is running on macOS, OpenRAG uses the `ocrmac` OCR engine. Other platforms use `easyocr`.

**Picture descriptions** adds image descriptions generated by the `SmolVLM-256M-Instruct` model to OCR processing. Enabling picture descriptions can slow ingestion performance.

## Use OpenRAG default ingestion instead of Docling serve

If you want to use OpenRAG's built-in pipeline instead of Docling serve, set `DISABLE_INGEST_WITH_LANGFLOW=true` in `Environment variables`.

The built-in pipeline still uses the Docling processor, but uses it directly without the Docling Serve API.

For more information, see `processors.py` in the OpenRAG repository.


## Knowledge ingestion flows

`Flows` in Langflow are functional representations of application workflows, with multiple `component` nodes connected as single steps in a workflow.

The **OpenSearch Ingestion** flow is the default knowledge ingestion flow in OpenRAG: when you **Add Knowledge** in OpenRAG, you run the OpenSearch Ingestion flow in the background. The flow ingests documents using **Docling Serve** to import and process documents.

This flow contains ten components connected together to process and store documents in your knowledge base.

- The **Docling Serve component** processes input documents by connecting to your instance of Docling Serve.
- The **Export DoclingDocument component** exports the processed DoclingDocument to markdown format with image export mode set to placeholder. This conversion makes the structured document data into a standardized format for further processing.
- Three **DataFrame Operations components** sequentially add metadata columns to the document data of `filename`, `file_size`, and `mimetype`.
- The **Split Text component** splits the processed text into chunks with a chunk size of 1000 characters and an overlap of 200 characters.
- Four **Secret Input** components provide secure access to configuration variables: `CONNECTOR_TYPE`, `OWNER`, `OWNER_EMAIL`, and `OWNER_NAME`. These are runtime variables populated from OAuth login.
- The **Create Data** component combines the secret inputs into a structured data object that will be associated with the document embeddings.
- The **Embedding Model component** generates vector embeddings using OpenAI's `text-embedding-3-small` model. The embedding model is selected at [Application onboarding] and cannot be changed.
- The **OpenSearch component** stores the processed documents and their embeddings in the `documents` index at `https://opensearch:9200`. By default, the component is authenticated with a JWT token, but you can also select `basic` auth mode, and enter your OpenSearch admin username and password.

All flows included with OpenRAG are designed to be modular, performant, and provider-agnostic. To modify a flow, click  **Settings**, and click **Edit in Langflow**. OpenRAG's visual editor is based on the [Langflow visual editor](#), so you can edit your flows to match your specific use case.

## OpenSearch URL Ingestion flow

An additional knowledge ingestion flow is included in OpenRAG, where it is used as an MCP tool by the **Open Search Agent flow**. The agent calls this component to fetch web content, and the results are ingested into OpenSearch.

For more on using MCP clients in Langflow, see [MCP clients](#).

To connect additional MCP servers to the MCP client, see [Connect to MCP servers from your application](#).

# Environment variables

OpenRAG recognizes [supported environment variables](#) from the following sources:

- [Environment variables](#) - Values set in the `.env` file.
- [Langflow runtime overrides](#) - Langflow components may tweak environment variables at runtime.
- [Default or fallback values](#) - These values are default or fallback values if OpenRAG doesn't find a value.

## Configure environment variables

Environment variables are set in a `.env` file in the root of your OpenRAG project directory.

For an example `.env` file, see [.env.example](#) in the [OpenRAG repository](#).

The Docker Compose files are populated with values from your `.env`, so you don't need to edit the Docker Compose files manually.

Environment variables always take precedence over other variables.

## Set environment variables

To set environment variables, do the following.

1. Stop OpenRAG.
2. Set the values in the `.env` file:

```
LOG_LEVEL=DEBUG
LOG_FORMAT=json
SERVICE_NAME=openrag-dev
```

3. Start OpenRAG.

Updating provider API keys or provider endpoints in the `.env` file will not take effect after [Application onboarding](#). To change these values, you must:

1. Stop OpenRAG.
  2. Remove the containers:
-

```
docker-compose down
```

3. Update the values in your `.env` file.
4. Start OpenRAG containers.

```
docker-compose up -d
```

5. Complete [Application onboarding](#) again.

## Supported environment variables

All OpenRAG configuration can be controlled through environment variables.

### AI provider settings

Configure which AI models and providers OpenRAG uses for language processing and embeddings. For more information, see [Application onboarding](#).

Variable	Default	Description
<code>EMBEDDING_MODEL</code>	<code>text-embedding-3-small</code>	Embedding model for vector search.
<code>LLM_MODEL</code>	<code>gpt-4o-mini</code>	Language model for the chat agent.
<code>MODEL_PROVIDER</code>	<code>openai</code>	Model provider, such as OpenAI or IBM watsonx.ai.
<code>OPENAI_API_KEY</code>	-	Your OpenAI API key. Required.
<code>PROVIDER_API_KEY</code>	-	API key for the model provider.
<code>PROVIDER_ENDPOINT</code>	-	Custom provider endpoint. Only used for IBM or Ollama providers.
<code>PROVIDER_PROJECT_ID</code>	-	Project ID for providers. Only required for the IBM watsonx.ai provider.

### Document processing

Control how OpenRAG processes and ingests documents into your knowledge base. For more information, see [Ingestion](#).

Variable	Default	Description
<code>CHUNK_OVERLAP</code>	<code>200</code>	Overlap between chunks.
<code>CHUNK_SIZE</code>	<code>1000</code>	Text chunk size for document processing.
<code>DISABLE_INGEST_WITH_LANGFLOW</code>	<code>false</code>	Disable Langflow ingestion pipeline.
<code>DOCLING_OCR_ENGINE</code>	-	OCR engine for document processing.
<code>OCR_ENABLED</code>	<code>false</code>	Enable OCR for image processing.
<code>OPENRAG_DOCUMENTS_PATHS</code>	<code>./documents</code>	Document paths for ingestion.
<code>PICTURE_DESCRIPTIONS_ENABLED</code>	<code>false</code>	Enable picture descriptions.

## Langflow settings

Configure Langflow authentication.

Variable	Default	Description
<code>LANGFLOW_AUTO_LOGIN</code>	<code>False</code>	Enable auto-login for Langflow.
<code>LANGFLOW_CHAT_FLOW_ID</code>	pre-filled	This value is pre-filled. The default value is found in <a href="#">.env.example</a> .
<code>LANGFLOW_ENABLE_SUPERUSER_CLI</code>	<code>False</code>	Enable superuser CLI.

Variable	Default	Description
<code>LANGFLOW_INGEST_FLOW_ID</code>	pre-filled	This value is pre-filled. The default value is found in <a href="#">.env.example</a> .
<code>LANGFLOW_KEY</code>	auto-generated	Explicit Langflow API key.
<code>LANGFLOW_NEW_USER_IS_ACTIVE</code>	False	New users are active by default.
<code>LANGFLOW_PUBLIC_URL</code>	<code>http://localhost:7860</code>	Public URL for Langflow.
<code>LANGFLOW_SECRET_KEY</code>	-	Secret key for Langflow internal operations.
<code>LANGFLOW_SUPERUSER</code>	-	Langflow admin username. Required.
<code>LANGFLOW_SUPERUSER_PASSWORD</code>	-	Langflow admin password. Required.
<code>LANGFLOW_URL</code>	<code>http://localhost:7860</code>	Langflow URL.
<code>NUDGES_FLOW_ID</code>	pre-filled	This value is pre-filled. The default value is found in <a href="#">.env.example</a> .



Variable	Default	Description
<code>SYSTEM_PROMPT</code>	"You are a helpful AI assistant with access to a knowledge base. Answer questions based on the provided context."	System prompt for the Langflow agent.

## OAuth provider settings

Configure OAuth providers and external service integrations.

Variable	Default	Description
<code>AWS_ACCESS_KEY_ID</code> / <code>AWS_SECRET_ACCESS_KEY</code>	-	AWS integrations.
<code>GOOGLE_OAUTH_CLIENT_ID</code> / <code>GOOGLE_OAUTH_CLIENT_SECRET</code>	-	Google OAuth authentication.
<code>MICROSOFT_GRAPH_OAUTH_CLIENT_ID</code> / <code>MICROSOFT_GRAPH_OAUTH_CLIENT_SECRET</code>	-	Microsoft OAuth.
<code>WEBHOOK_BASE_URL</code>	-	Base URL for webhook endpoints.

## OpenSearch settings

Configure OpenSearch database authentication.

Variable	Default	Description
<code>OPENSEARCH_HOST</code>	<code>localhost</code>	OpenSearch host.
<code>OPENSEARCH_PASSWORD</code>	-	Password for OpenSearch admin user. Required.
<code>OPENSEARCH_PORT</code>	<code>9200</code>	OpenSearch port.
<code>OPENSEARCH_USERNAME</code>	<code>admin</code>	OpenSearch username.

## System settings

Configure general system components, session management, and logging.

Variable	Default	Description
LANGFLOW_KEY_RETRIES	15	Number of retries for Langflow key generation.
LANGFLOW_KEY_RETRY_DELAY	2.0	Delay between retries in seconds.
LANGFLOW_VERSION	latest	Langflow Docker image version.
LOG_FORMAT	-	Log format (set to "json" for JSON output).
LOG_LEVEL	INFO	Logging level (DEBUG, INFO, WARNING, ERROR).
MAX_WORKERS	-	Maximum number of workers for document processing.
OPENRAG_VERSION	latest	OpenRAG Docker image version.
SERVICE_NAME	openrag	Service name for logging.
SESSION_SECRET	auto-generated	Session management.

## Langflow runtime overrides

Langflow runtime overrides allow you to modify component settings at runtime without changing the base configuration.

Runtime overrides are implemented through **tweaks** - parameter modifications that are passed to specific Langflow components during flow execution.

For more information on tweaks, see [Input schema \(tweaks\)](#).

## Default values and fallbacks

When no environment variables or configuration file values are provided, OpenRAG uses default values. These values can be found in the code base at the following locations.

## OpenRAG configuration defaults

These values are defined in `config_manager.py` in the OpenRAG repository.

## **System configuration defaults**

These fallback values are defined in `settings.py` in the OpenRAG repository.

# Troubleshoot

This page provides troubleshooting advice for issues you might encounter when using OpenRAG or contributing to OpenRAG.

## OpenSearch fails to start

Check that `OPENSEARCH_PASSWORD` set in [Environment variables](#) meets requirements. The password must contain at least 8 characters, and must contain at least one uppercase letter, one lowercase letter, one digit, and one special character that is strong.

## Langflow connection issues

Verify the `LANGFLOW_SUPERUSER` credentials set in [Environment variables](#) are correct.

## Memory errors

### Container out of memory errors

Increase Docker memory allocation or use [docker-compose-cpu.yml](#) to deploy OpenRAG.

### Podman on macOS memory issues

If you're using Podman on macOS, you may need to increase VM memory on your Podman machine. This example increases the machine size to 8 GB of RAM, which should be sufficient to run OpenRAG.

```
podman machine stop
podman machine rm
podman machine init --memory 8192 # 8 GB example
podman machine start
```

## Port conflicts

Ensure ports 3000, 7860, 8000, 9200, 5601 are available.

## Langflow container already exists

If you are running other versions of Langflow containers on your machine, you may encounter an issue where Docker or Podman thinks Langflow is already up.

Remove just the problem container, or clean up all containers and start fresh.

To reset your local containers and pull new images, do the following:

1. Stop your containers and completely remove them.

**Podman**     **Docker**

```
# Stop all running containers
docker stop $(docker ps -q)
# Remove all containers (including stopped ones)
docker rm --force $(docker ps -aq)
# Remove all images
docker rmi --force $(docker images -q)
# Remove all volumes
docker volume prune --force
# Remove all networks (except default)
docker network prune --force
# Clean up any leftover data
docker system prune --all --force --volumes
```

2. Restart OpenRAG and upgrade to get the latest images for your containers.

```
uv sync
uv run openrag
```

3. In the OpenRAG TUI, click **Status**, and then click **Upgrade**. When the **Close** button is active, the upgrade is complete. Close the window and open the OpenRAG application.