



# Vectoria: Supercomputer-driven vector database for private LLM retrieval

User Manual, version 1.0



**EuroHPC**  
Joint Undertaking

Funded by the European Union. This work has received funding from the European High Performance Computing Joint Undertaking (JU) and Germany, Bulgaria, Austria, Croatia, Cyprus, Czech Republic, Denmark, Estonia, Finland, Greece, Hungary, Ireland, Italy, Lithuania, Latvia, Poland, Portugal, Romania, Slovenia, Spain, Sweden, France, Netherlands, Belgium, Luxembourg, Slovakia, Norway, Türkiye, Republic of North Macedonia, Iceland, Montenegro, Serbia under grant agreement No 101101903.

This page intentionally blank

## Summary

Acronyms and Terminology.....	4
1. Introduction.....	5
2. EuroCC: Objectives.....	5
3. Service Description .....	6
4. Architectural Diagram .....	6
Tasks .....	6
4.1.1 Build vector db.....	7
4.1.2 Inference.....	7
5. Self-Deployment Guide .....	9
5.1 Prerequisites.....	9
5.2 Step-by-step Installation & Deploy .....	9
5.3 USAGE.....	10
5.3.1 Task Configuration .....	10
Application information and contributions .....	14

## Table of Figures

FIGURE 1: NCCs ACROSS EUROPE.....	5
FIGURE 2: BUILD VECTOR DB WORKFLOW .....	7
FIGURE 3: INFERENCE WORKFLOW.....	7

## Acronyms and Terminology

AI	Artificial Intelligence	Field of study in computer science that develops and studies intelligent machines.
LLM	Large Language Model	A type of AI model trained on vast amounts of text to perform natural language understanding and generation.
RAG	Retrieval-Augmented Generation	A technique that combines information retrieval and generative AI for producing contextually relevant outputs.
CLI	Command-Line Interface	A text-based interface used to interact with software or operating systems.
FAISS	Facebook AI Similarity Search	A library for fast nearest-neighbor search and vector similarity.
API	Application Programming Interface	A set of rules and tools for building and interacting with software applications.
GPU	Graphics Processing Unit	Specialized hardware designed for accelerating computational tasks, especially in AI and ML.
MMR	Maximal Marginal Relevance	A ranking technique that optimizes information relevance and diversity in retrieved results.
RAGAS	Retrieval-Augmented Generation Assessment	A tool or metric for evaluating the performance of RAG-based systems.
HPC	High-Performance Computing	Computing systems with high-speed processing capabilities used for complex computations.
Embedding	-	A vector representation of data (e.g., text) used for similarity searches and machine learning tasks.
Inference	-	The process of generating predictions or outputs from a trained model.

## 1. Introduction

Within the framework of the European Project EuroCC2, the national competence center EuroCC Italy developed “Vectoria: Supercomputer-driven vector database for private LLM retrieval”. This project aims to equip small and medium enterprises with a service to streamline access to private, internal documentation through intuitive chatbot-like interactions. This initiative targets the challenges of navigating unstructured and dispersed knowledge stored across diverse organizational repositories, offering an efficient solution for users seeking concise, contextually relevant answers.

By employing state-of-the-art retrieval and generative AI technologies, the system ensures privacy preservation while delivering high-quality responses. This document aims to outline the requirements, core functionalities, and operational guidance for deploying and utilizing the system effectively.

The service is designed to integrate with High Performance Computing (HPC) systems, leveraging the computational power of supercomputers to enhance the speed and precision of a generative model inference. This approach enables scalable, private, and high-performance information retrieval, tailored to meet the needs of users in a data-rich environment.

## 2. EuroCC: Objectives

[EuroCC 2](#) is a project which works to identify and address the skills gaps in the European High-Performance Computing (HPC) ecosystem and coordinate cooperation across Europe to ensure a consistent skills base. The mission of EuroCC is to improve the technological readiness of Europe. In particular, the role of EuroCC2 is to establish and run a network of more than 30 NCCs across the EuroHPC Participating States. The NCCs act as single points of access in each country between stakeholders and national and EuroHPC systems. They operate on a regional and national level to liaise with local communities, in particular SMEs, map HPC competencies and facilitate access to European HPC resources for users from the private and public sector.



Figure 1: NCCs across Europe

EuroCC2 delivers training, interacts with industry, develops competence mapping and communication materials and activities, and supports the adoption of HPC services in other related fields, such as quantum computing, artificial intelligence (AI), high performance data analytics (HPDA) to expand the HPC user base.

### 3. Service Description

The “Vectoria: Supercomputer-driven vector database for private LLM retrieval” project introduces a tailored approach to optimize knowledge access within organizations. By focusing on the integration of advanced AI-driven methods, it enables users to interact seamlessly with internal documentation via chatbot-like interfaces. The project emphasizes creating a unified, efficient framework to navigate unstructured data by delivering concise, contextually relevant answers to users while upholding stringent privacy standards.

#### Key Objectives:

- **Efficient Knowledge Management:** scattered internal documentation is unified into a knowledge source accessible through querying.
- **High Precision and Relevance:** Deliver results that precisely match user queries by using vector-based retrieval techniques.
- **Data Privacy:** Operate entirely within the organization's infrastructure to ensure compliance with internal security protocols.

### 4. Architectural Diagram

The RAG pipeline consists of several tightly integrated components, each playing a critical role in ensuring the system's accuracy. Below is a comprehensive breakdown of its architecture and workflow.

#### System Overview

The architecture comprises the following core modules:

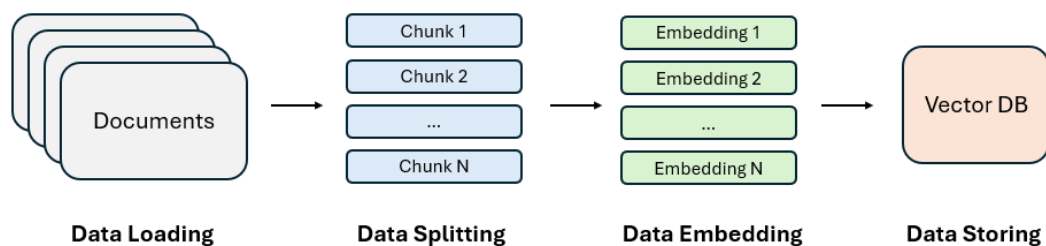
1. **Configuration:** The system can be easily customized through a yaml-based configuration file.
2. **Document Preprocessor:** Handles ingestion, cleaning, and splitting of internal documents into smaller portions of text (*chunks*) and transforming these chunks into embeddings.
3. **Vector Database:** Stores vector embeddings of the document chunks for efficient similarity searches.
4. **Retriever:** Locates relevant chunks from the vector database based on user queries.
5. **Post-retrieval strategies:** Includes re-ranker and full paragraph retriever.
6. **Generative Model:** Synthesizes human-like answers using a pre-trained large language model (e.g., GPT).

#### Tasks

The RAG pipeline provides two primary operational tasks:

- **Build vector database:** Responsible for preparing and indexing the documents to create a vector database.
- **Inference:** Handles user queries by retrieving and generating context-aware responses.

### 4.1.1 Build vector db



This task focuses on the ingestion and indexing of internal documents. It transforms a set of documents into a searchable structure, ensuring that subsequent retrieval operations are both fast and accurate. In technical language (and inside the codebase) this operation is also referred as “Build Index”.

#### Workflow:

Figure 2: Build vector db workflow

**Prerequisites:** as version 1.0, supported formats include PDF and DOCX.

1. **Data Preprocessing (Loading and Splitting):** Clean the text to remove irrelevant elements like boilerplate text, extra whitespace, or special characters. Split documents into chunks of manageable size (e.g., 512 characters) to optimize retrieval granularity. Parameters for cleaning and splitting can be configured using a single configuration file.
2. **Data Embedding:** Generate vector embeddings for each chunk using pre-trained models (e.g., bge-m3). Each embedding is a n-dimensional vector which captures the semantic meaning of the chunk.
3. **Data Storing:** Store the embeddings in a vector database (e.g., FAISS vector database) dumped on disk. Attach metadata (e.g., source file, paragraph number) for efficient filtering and explainability.

### 4.1.2 Inference

The inference task enables the chatbot to provide accurate answers by combining retrieval and generation. It works exploiting the indexed data to ensure high-quality responses.

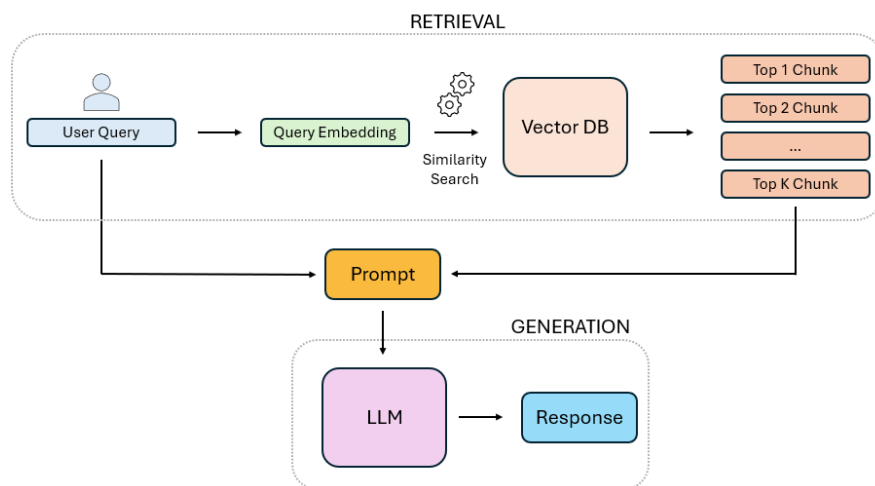


Figure 3: Inference workflow

**Workflow:**

1. **Query Reception:** User submits a question via a CLI or file
2. **Query Embedding:** Convert the user query into a vector embedding using the same pre-trained model used during indexing.
3. **Vector Search:** Search the vector database for the most semantically similar chunks to the query. Retrieve the top-k chunks (e.g., top 5 most relevant) using a similarity metric (e.g., cosine similarity).
4. **Post-Retrieval strategies:** Perform re-ranking or full paragraph retrieving, if necessary.
5. **Prompt Construction:** Combine the *retrieved chunks* with the *question* and the *instructions* to form a context window for the generative model.
6. **Answer Generation:** Pass the query and context to a generative AI model (e.g., llama3.1) to generate a natural language response tailored to answer the user's query.



## 5. Self-Deployment Guide

The following provides a detailed description of how to use the script provided to deploy the service.

### 5.1 Prerequisites

The following prerequisites must be met before continuing:

- System requirements:
  - Minimum: **Embedder** (bge-m3) + **LLM** (Llama-3.1-8B-Instruct)
    - **VRAM**: ~36GB
  - Recommended: **Embedder** (bge-m3) + **LLM** (Llama-3.1-70B-Instruct)
    - **VRAM**: ~252GB
- The user setting up the service must have a Linux terminal with Python 3.8 or later.
- If installing on a remote host, SSH access is required.
- For some LLMs, you will also need to register your SSH key to access the download. Consult the documentation of your LLM of choice for instructions.
- `git-lfs` is required to download some larger LLMs. Make sure it is installed and available on the install destination.

### 5.2 Step-by-step Installation & Deploy

To install Vectoria follow the steps described below.

1. Clone the repository containing the Ansible codebase:

```
git clone https://github.com/Eurocc-Italy/Vectoria.git
```

2. Create and activate a Python virtual environment:

```
python3 -m venv virtualenv
source virtualenv/bin/activate
```

3. Move to the Vectoria ansible folder and install the necessary requirements:

```
cd vectoria/ansible
pip install -r requirements.txt
```

4. In the `inventory/hosts` file, edit the `ansible_host` and `ansible_user` variables with the hostname (or IP address) and the username of the install destination to be used, respectively.
5. If installing locally, edit the `deploy.yml` file, changing the following `hosts` variable from `hpc` to `localhost`.
6. Edit the `inventory/vars/config.yml` file as necessary. Comments explaining what each variable pertains to are provided in the file itself. Some reasonable defaults are provided. For the inference engine of choice, we recommend the following models based on the scenario:
  - Consumer PC (~8GB vRAM): Llama-3.2-1B-Instruct
  - High-end workstation (~32 GB vRAM): Llama-3.1-8B-Instruct
  - HPC cluster (~256 GB vRAM): Llama-3.1-70B-Instruct
7. Run the Ansible playbook:

```
ansible-playbook -i inventory deploy.yml
```

- Alternatively, it is also possible to install only the Vectoria library locally via pip (editable install mode is required), and configure each setting manually:

```
pip install -e vectoria
```

## 5.3 USAGE

In this section examples are provided in order to launch tasks.

### 1. Build Index:

```
python vectoria \
# Launch vectoria command
--config 'etc/default/default_config.yaml' \
# [Required: TRUE] Specify config file PATH
build_index \
# [Required: TRUE] Specify task to be executed
--input-docs-dir 'test/data' \
# [Required: TRUE] Specify dir containing input files to be indexed
--output-dir 'test/index' \
# [Required: TRUE] Specify dir that will contain the new index
--output-suffix '_my_test_index'
# [Required: FALSE] Specify suffix to add to the output file name to customize the index
```

### 2. Inference:

```
python vectoria \
# Launch vectoria command
--config 'etc/default/default_config.yaml' \
# [Required: TRUE] Specify config file PATH
inference \
# [Required: TRUE] Specify task to be executed
--index-path 'test/data/index/my_new_test_index' \
# [Required: TRUE] Specify pickle file PATH to an already existing index
--test-set-path 'test/data/results/my_questions.json' \
# [Required: FALSE] Specify questions file PATH
--questions 'Who is the CEO?' 'Who is the VP?' \
# [Required: FALSE] Specify questions directly via CLI
--output-dir 'test/data/results'
# [Required: TRUE] Specify dir that will contain a file with the answers
```

N.B.: `--test-set-path` and `--questions` are both not formally *strictly* required, since vectoria is able to handle a dual mode. Using `--questions` argument will override questions from `--test-set-path` and will provide answers only through CLI, avoiding any dump on file. IMPORTANT: at least one MUST be provided.

### 5.3.1 Task Configuration

This section explains a configuration file example. It contains every parameter that can be customized to adapt the behavior of the system to the desired one, covering several components in detail.

PARAMETERS	TYPE	VALUES (default in bold)	DESCRIPTION
<code>vectoria_logs_dir:</code>	str	<code>vectoria_logs</code>	PATH where log files will be dumped
<code>log_level:</code>	str	<code>[DEBUG   INFO   CRITICAL]</code>	Sets general log level
<code>langchain_tracking:</code>	str	<code>[true   false]</code>	Enables LangChain tracking backend
<code>system_prompts_lang:</code>	str	<code>[eng   it]</code>	Language for system prompts, e.g., English or Italian
<code>data_ingestion:</code>			
<code>multiprocessing:</code>	str	<code>[true   false]</code>	Enables/disables parallel processing for ingestion tasks
<code>extraction:</code>			
<code>format:</code>	str	<code>[docx   pdf]</code>	Specifies input document format
<code>dump_doc_structure_on_file:</code>	str	<code>[true   false]</code>	Saves document structure to file if true
<code>regexes_for_metadata_extraction:</code>			List of regex patterns for metadata extraction
<code>- name:</code>	str	<code>DOC_ID</code>	Name of the metadata field
<code>pattern:</code>	str	<code>'^Document Title'</code>	Regex pattern to match the field
<code>regexes_for_replacement:</code>			List of regex rules for text preprocessing
<code>- name:</code>	str	<code>remove_multiple_spaces</code>	Rule name
<code>pattern:</code>	str	<code>'[ \t]{2,}'</code>	Regex to find multiple spaces/tabs
<code>replace_with:</code>		<code>' '</code>	Replacement value
<code>- name:</code>	str	<code>remove_bullets</code>	Rule name
<code>pattern:</code>	str	<code>'^\s*[\u2022\u25AA\u27A2]\s*'</code>	Regex to match bullets
<code>replace_with:</code>		<code>''</code>	Replace bullets with empty string
<code>- name:</code>	str	<code>remove_ligature_st</code>	Rule name
<code>pattern:</code>	str	<code>'st'</code>	Regex to match 'st'
<code>replace_with:</code>	str	<code>'st'</code>	Replace 'st' with 'st'
<code>chunking:</code>			
<code>chunk_size:</code>	int	<code>[256   512   1024]</code>	Size of each chunk for processing
<code>chunk_overlap:</code>	int	<code>[128   256   512]</code>	Overlap between consecutive chunks
<code>separators:</code>	List[str]	<code>["\n\n", "\n", " ", ""]</code>	List of separators for chunking
<code>is_separator_regex:</code>	List[str]	<code>[false, false, false, false]</code>	Whether separators are regex patterns
<code>dump_chunks_on_file:</code>	str	<code>[true   false]</code>	Save chunked output to file if true

<b>vector_store:</b>			
name:	str	faiss	Vector store backend (e.g., FAISS)
model_name:	str	BAAI/bge-m3	Model used for generating embeddings
device:	str	[cuda   cpu]	Device for vector store operations
normalize_embeddings:	str	false	Normalizes embeddings if true
<b>retriever:</b>			
enabled:	str	[true   false]	Enables/disables the retriever
top_k:	int	5	Number of top results to retrieve
search_type:	str	'mmr'	Search type (e.g., Maximal Marginal Relevance)
fetch_k:	int	5	Number of documents to fetch before filtering
lambda_mult:	float	0.5	Lambda parameter for MMR search
<b>reranker:</b>			
enabled:	str	[true   false]	Enables/disables reranking
reranked_top_k:	int	3	Number of reranked results to return
<b>inference_engine:</b>			Inference engine configuration for reranking
name:	str	[huggingface   ollama   openai   vllm]	Name of inference engine backend
url:	str	null	URL for the inference API endpoint (e.g. http://localhost:8898/v1)
api_key:	str	null	API key for the inference engine
model_name:	str	BAAI/bge-reranker-v2-gemma	Model name for reranking
device:	str	[cuda   cpu]	Device for inference engine
load_in_4bit:	str	[true   false]	Enables 4-bit model loading for optimization
load_in_8bit:	str	[true   false]	Enables 8-bit model loading for optimization
max_new_tokens:	int	150	Maximum tokens generated in output
trust_remote_code:	str	[true   false]	Trust code from remote repositories if true
device_map:	str	[auto   null]	Device mapping for distributed inference
temperature:	float	0.1	Sampling temperature for generation
<b>full_paragraphs_retriever:</b>			
enabled:	str	[true   false]	Enables/disables retrieval of full paragraphs
<b>inference_engine:</b>			General inference engine configuration

<b>name:</b>	str	[huggingface   ollama   openai   vllm]	Name of the inference engine
<b>url:</b>	str	null	URL for the inference engine API
<b>api_key:</b>	str	null	API key for the inference engine
<b>model_name:</b>	str	meta-llama/Meta-Llama-3.1-8B-Instruct	Model name for inference
<b>device:</b>	str	[cuda   cpu]	Device for inference
<b>load_in_8bit:</b>	str	[true   false]	Enables 8-bit model loading for optimization
<b>max_new_tokens:</b>	int	150	Maximum tokens generated in output
<b>trust_remote_code:</b>	str	[true   false]	Trust code from remote repositories if true
<b>device_map:</b>	str	[auto   null]	Automatic device mapping for distributed inference
<b>temperature:</b>	float	0.1	Sampling temperature for generation
<b>evaluation:</b>			Evaluation configuration
<b>tool:</b>	str	ragas	Evaluation tool name
<b>inference_engine:</b>			Inference engine for evaluation
<b>name:</b>	str	vllm	Name of the inference engine backend
<b>model_name:</b>	str	hugging-quants/Meta-Llama-3.1-8B-Instruct-AWQ-INT4	Model for evaluation
<b>url:</b>	str	null	API endpoint for inference engine
<b>api_key:</b>		null	API key for the inference engine
<b>embeddings_engine:</b>			Embedding engine for evaluation
<b>name:</b>	str	vllm	Name of the embedding engine
<b>model_name:</b>	str	BAAI/bge-multilingual-gemma2	Model for embedding generation
<b>url:</b>	str	null	API endpoint for embeddings engine
<b>api_key:</b>	str	null	API key for the embeddings engine

## Application information and contributions

The application and this manual are released under the MIT open source license.

Copyright (c) 2024 EuroCC Italy, Cineca, Leonardo S.p.A.

Permission is hereby granted, free of charge, to any person obtaining a copy of this software and associated documentation files (the "Software"), to deal in the Software without restriction, including without limitation the rights to use, copy, modify, merge, publish, distribute, sublicense, and/or sell copies of the Software, and to permit persons to whom the Software is furnished to do so, subject to the following conditions:

The above copyright notice and this permission notice shall be included in all copies or substantial portions of the Software.

THE SOFTWARE IS PROVIDED "AS IS", WITHOUT WARRANTY OF ANY KIND, EXPRESS OR IMPLIED, INCLUDING BUT NOT LIMITED TO THE WARRANTIES OF MERCHANTABILITY, FITNESS FOR A PARTICULAR PURPOSE AND NONINFRINGEMENT. IN NO EVENT SHALL THE AUTHORS OR COPYRIGHT HOLDERS BE LIABLE FOR ANY CLAIM, DAMAGES OR OTHER LIABILITY, WHETHER IN AN ACTION OF CONTRACT, TORT OR OTHERWISE, ARISING FROM, OUT OF OR IN CONNECTION WITH THE SOFTWARE OR THE USE OR OTHER DEALINGS IN THE SOFTWARE.

This manual, *version 1.0*, refers to tag *v1.0* of repos:

<https://github.com/Eurocc-Italy/Vectoria>



The **Vectoria** application was developed by **EuroCC Italy** national competence center, with the specific contribution of **CINECA** and **Leonardo S.p.A.**

The project was realized thanks to:

Luca Babetto, Leonardo Baroncelli, Carolina Berucci, Chiara Malizia, Eric Pascolo, Andrea Proia

Thanks to SAVIA Team for the technical input (<https://assemblealegislativa.github.io/savia/>)