# Environment and Tooling Milestone

See the project README.md file for all instructions for running the model and how it works. Consider the README.md a more detailed report. This project report only shows screenshots proving each milestone is complete.

The screenshot below shows the containers started by the project docker compose file(mongo, gradio, qdrant, ollama). ClearML is run from the web instead as shown in the later milestones



After running docker compose up, you can see the containers respond whenever running files in the project folder if setup correctly. (Ignore the output jargon)

# ETL Milestone

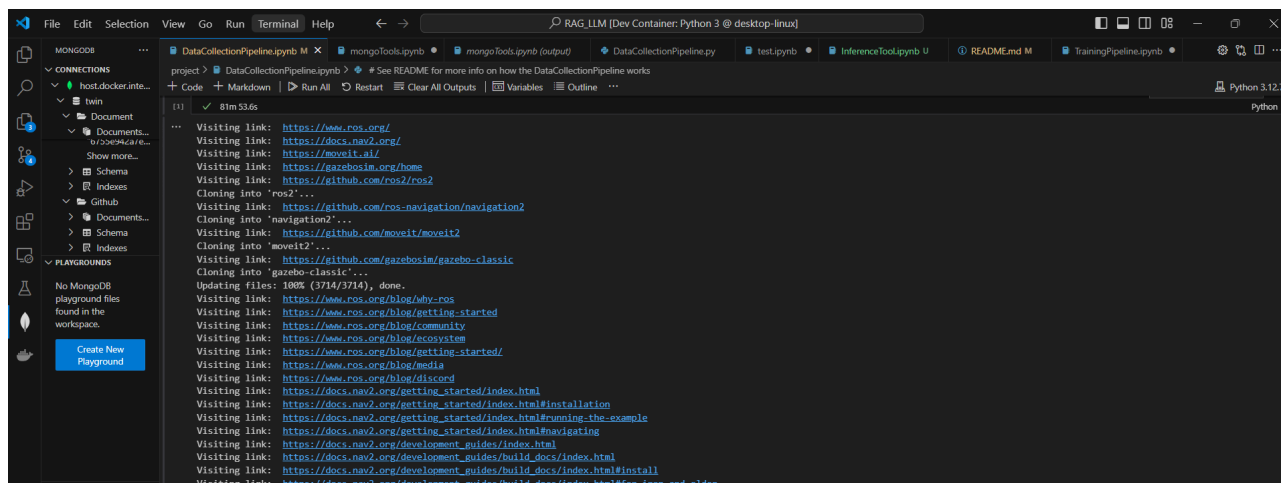The image below shows the ClearML/DataCollectionPipeline.py result in the ClearML website. It takes in links as parameters and outputs a list of texts from documents(documentation) and codes(github files) that come from the links it crawled. ClearML was run only on files related to the "ros" domain to save time. You can run either the ClearML or the ipynb notebook on the full input since they do the same thing, but ClearML requires setup and takes longer while the ipynb prints some output to the terminal, thus ipynb is preferred. Next I will go over the ipynb



The image below shows part of the output after running the DataCollectionPipeline.ipynb on the full list of links. It took an entire 81 minutes, so it might make sense to test the model on only the "ros" links(and ask it questions related to only that link). Related code shown next

Consider changing the links in the DataCollection.ipynb file to be similar to those in the ClearML/DataCollection.py file before running it to search only the ROS subdomain.



Alternatively, you can replace the following lines with "pass" to avoid searching all links inside each website(which each would have links of their own and so on, creating a massive tree of links). This large tree of links formed is the main reason the crawling takes over an hour.



```
#print("Adding subdirectory: ", link + subdirectory)
links.append(newLink)
```

The image below shows Tools/mongoTools.ipynb giving a summary of what was stored in the mongoDB. In total, 453 documents and 6800 code files were processed. And samples of what the documents and code looks like is shown at the bottom. Each piece of data contains an id, link, document type, and text content.

# Featurization Pipeline Milestone

The ClearML/FeaturePipeline.py file retrieves documents, cleans them, chunks them, creates embeddings, then stores those embeddings in Qdrant as shown in the ClearML website below



The image below shows the FeaturePipeline.ipynb being run. The pipeline took over 2 hours to run(depends on links crawled from DataPipeline). It took so long, the Qdrant cursor/iterator timed out, skipping some links.

The image below shows the first cell of Tools/qdrantTools.ipynb being used to show summary of the Qdrant vector database after running the FeaturePipeline.ipynb notebook to store embeddings in it.



The image below shows the second cell of Tools/qdrantTools.ipynb being used to show how many embeddings were completed.



The third cell shows a proof of concept for cosine distance(used to retrieve embeddings)

# Deploying the App Milestone

To generate an answer: the ClearML/InferencePipeline.py file performs query expansion to generate a related query (rewording query may lead to retrieving more useful embeddings), performs self-querying to get metadata about the query, searches for related vectors/embeddings compared to the query and related query(from query expansion), concatenates the results of the searches, reranks the results to filter out the most useful ones, builds a prompt using the query metadata and results, then feeds the prompt to the model to return the answer. The image below shows the ClearML website giving a overview of what this process looks like



The image below shows Tools/InferenceTool.ipynb showing the output of each step that was used to generate an answer to a sample prompt along with a comparison to the original model. You can change the query to whatever you want and it will give a similar answer as gradio. The first link confirms that the answer is correct while the baseline model did not know the answer!

The image below shows the gradio app being run and its response to a sample query. Additionally, you can see that the docker containers from docker compose up responded when the query was submitted. Alternatively, you can launch app.py, or even run the Tools/InferenceTools.ipynb file with a custom query. There are many ways to get an output as long as ollama is being served either by docker or in your terminal.



The images below show the model's responses to various queries

## Screenshot 1

Sample Prompt

How can I develop the navigation stack of an agent with egomotion? ▼

☑ useSample

Query

[ ]

| Clear | Submit |

output

Developing the navigation stack of an agent with egomotion involves integrating a combination of sensors, mapping algorithms, and motion planning techniques. Egomotion provides the capability to estimate the robot's pose and orientation using cameras or other vision systems, which is essential for building a robust navigation system.

Here are some steps you can take to develop the navigation stack:

1. **Data Collection**: Gather data on the environment, including the layout, obstacles, and any notable features. This data will be used to build and update the mapping system.

2. **Mapping Algorithm**: Implement a mapping algorithm that can create an accurate representation of the environment. Some popular algorithms include SLAM (Simultaneous Localization and Mapping), Monocular-SLAM, and GraphSLAM.

3. **Sensor Integration**: Integrate the egomotion sensor data with other sensors such as lidar or radar to improve navigation accuracy.

4. **Motion Planning**: Implement a motion planning algorithm that can generate efficient paths for the robot. Some popular algorithms include PRM (Probabilistic Roadmap), Rapidly-exploring Random Tree (RRT), and Motion Planning under Uncertainty (MPU).

Flag

## Screenshot 2

Sample Prompt

How can I develop the navigation stack of an agent with egomotion? ▼

☑ useSample

Query

[ ]

| Clear | Submit |

output

4. **Motion Planning**: Implement a motion planning algorithm that can generate efficient paths for the robot. Some popular algorithms include PRM (Probabilistic Roadmap), Rapidly-exploring Random Tree (RRT), and Motion Planning under Uncertainty (MPU).

5. **Local Navigation**: Develop local navigation techniques, such as dead reckoning or wheel odometer estimation, to navigate within a small area.

6. **Global Navigation**: Implement global navigation techniques, such as SLAM or visual SLAM, to build a comprehensive map of the environment and enable long-range navigation.

7. **Optimization and Tuning**: Optimize and tune the navigation system for better performance, including adjusting parameters, fine-tuning algorithms, and incorporating feedback from the agent's experiences.

8. **Testing and Evaluation**: Test and evaluate the navigation system in different scenarios, such as varying environments, lighting conditions, or robot speeds, to ensure its robustness and reliability.

By following these steps, you can develop a comprehensive navigation stack for an agent with egomotion that enables effective local and global navigation in various environments.

Flag

## Window 1

Gradio | localhost:7860

**Sample Prompt**

What is ROS? ▼

☑ useSample

**Query**

[ ]

Clear | Submit

**output**

ROS stands for Robot Operating System. It's an open-source software framework that provides a common interface for software developers to create robot applications. ROS allows developers to write software that can run on any device with a C++ compiler, making it easy to implement robotics projects across different platforms.

ROS includes various components such as MoveIt (a motion planning library), Gazebo (an open-source simulation environment), and others, which enable robots to interact with their environment. It's widely used in the robotics community for building and controlling complex robotic systems.

The MoveIt package is particularly useful for executing time-parameterized joint trajectories to low-level hardware controllers through common interfaces. This allows robots to be controlled and programmed using a variety of planning algorithms and visualization tools, making it easier to develop and test robotics applications.

Flag

Use via API 🚀 · Built with Gradio 🟠

## Window 2

Gradio | localhost:7860

**Sample Prompt**

How many companies is Nav2 trusted by worldwide? ▼

☑ useSample

**Query**

[ ]

Clear | Submit

**output**

Nav2 is trusted by over 100 companies worldwide.

Flag

Use via API 🚀 · Built with Gradio 🟠

## Screenshot 1

Gradio — localhost:7860

**Sample Prompt**

How would I build a ROS 2 Navigation Framework and System? ▾

☑ useSample

**Query**

[                                        ]

[ Clear ]          [ Submit ]

**output**

Building a ROS 2 navigation framework and system involves several key components and steps. Here's an overview to get you started:

1. **Perception**: This is the process of collecting and interpreting data from sensors such as cameras, lidars, or radar. You can use pre-existing perception packages in ROS 2 like `nav2` or create your own using libraries like OpenCV or PCL.

2. **Planning**: After collecting sensor data, you need to plan a navigation path for your robot. ROS 2's `nav2` package provides a range of planning algorithms and tools to help with this step.

3. **Control**: Once you have a planned path, you need to control the robot's movement. This involves sending commands to the robot's actuators, which can be done using ROS 2's `ros_control` package.

4. **Localization**: In order for your navigation system to work effectively, it needs to know its location and orientation in the environment. You can use pre-existing localization packages like `nav2` or create your own using libraries like SLAM (Simultaneous Localization and Mapping) algorithms.

[ Flag ]

Use via API ⚡ · Built with Gradio 🟠

## Screenshot 2

Gradio — localhost:7860

**Sample Prompt**

How would I build a ROS 2 Navigation Framework and System? ▾

☑ useSample

**Query**

[                                        ]

[ Clear ]          [ Submit ]

**output**

5. **Visualization**: Visualizing the navigation process can help you understand what's happening and diagnose any issues that may arise. ROS 2 provides tools like RViz for this purpose.

6. **Integration with Physical Hardware**: Finally, you need to integrate your navigation system with physical hardware. This involves setting up the necessary actuators, sensors, and communication protocols to work seamlessly together.

To build a ROS 2 navigation framework and system, follow these general steps:

* Install ROS 2 on your computer.
* Choose a suitable navigation package (like `nav2`) and set it up according to its documentation.
* Implement perception, planning, control, localization, visualization, and integration with physical hardware components as needed.
* Test your navigation system thoroughly to ensure it's working correctly.

Keep in mind that building a full-fledged ROS 2 navigation framework and system can be complex and time-consuming. Be prepared to spend some time learning about ROS 2 and the specific packages you're using.

[ Flag ]

Use via API ⚡ · Built with Gradio 🟠

Sample Prompt

Write me code to move a robot using MoveIt

useSample

Query

Clear    Submit

output

To use MoveIt with your ROS environment, you'll need to have the necessary dependencies installed and configured. Here's an example of how you can achieve this:

```cpp
#include <moveit_commander/moveit_commander.h>
#include <moveit_planning_scene/planning_scene.h>
#include <move_group_interface/move_group_interface.h>

int main(int argc, char **argv) {
    // Initialize MoveIt
    moveit_commander::MoveGroupInterface* move_group = new moveit_commander::MoveGroupInterface("arm_group");

    // Set the robot's current pose
    geometry_msgs::Pose pose;
    pose.position.x = 0.5;  // X-coordinate of the pose
    pose.position.y = 0.5;  // Y-coordinate of the pose
    pose.position.z = 1.0;  // Z-coordinate of the pose

    move_group->setPose(pose);
```

Flag

Sample Prompt

Write me code to move a robot using MoveIt

useSample

Query

Clear    Submit

output

```cpp
    pose.position.y = 0.5;  // Y-coordinate of the pose
    pose.position.z = 1.0;  // Z-coordinate of the pose

    move_group->setPose(pose);

    // Wait for a short period to ensure the robot has finished moving
    ros::Duration(2.0).sleep();

    return 0;
}
```

This code snippet initializes MoveIt and sets the robot's current pose using `moveit_commander::MoveGroupInterface`. It then waits for a short period of time before exiting.

Remember that you need to replace "arm_group" with the actual name of your group. You can get this name from your MoveIt configuration file.

For more detailed instructions, please refer to [ROS documentation] (https://wiki.ros.org/moveit)
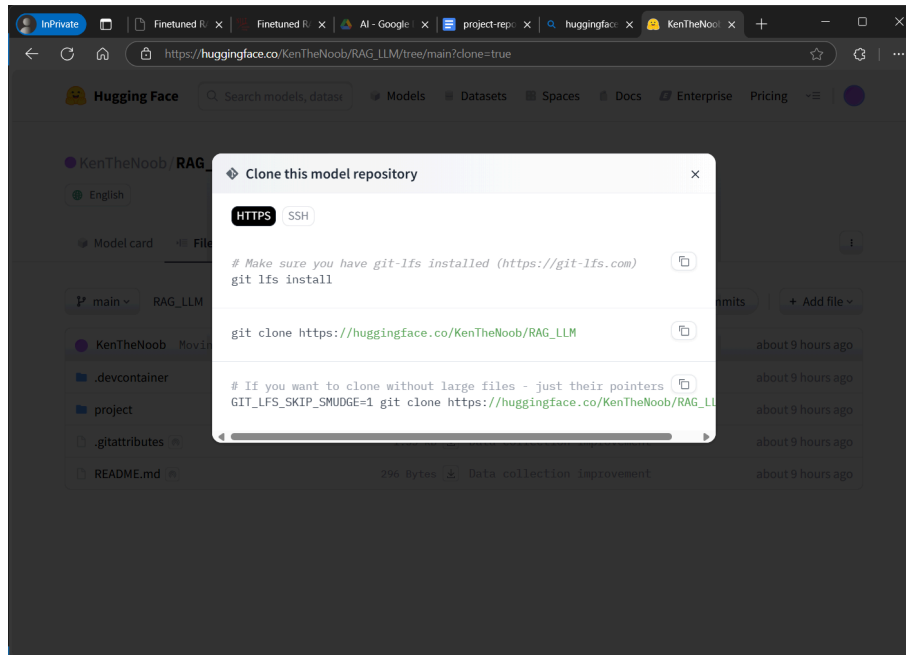
Flag

The image below shows the model uploaded and pullable on huggingface

Alternatively, you can pull from my github in the README.md, in which case, you start the devcontainer inside the project folder, not the base folder.