

Multimodal Visual Understand (Text Version)

Multimodal Visual Understand (Text Version)

- 1. Course Content
- 2. Preparation
 - 2.1 Content Description
- 3. Running the Example
 - 3.1 Starting the Program
 - 3.2 Test Examples
 - 3.2.1 Example 1
 - 3.2.2 Example 2
- 4. Source Code Analysis
 - action_service.py
 - model_service.py


1. Course Content


- 1. Learn to use the robot's visual understanding capabilities
- 2. Study new key source code

2. Preparation

2.1 Content Description

This course uses the Jetson Orin NANO as an example. For Raspberry Pi and Jetson Nano boards, you need to open a terminal on the computer and enter the command to enter the Docker container. Once inside the Docker container, enter the commands mentioned in this course in the terminal. For instructions on entering the Docker container from the computer, refer to **[01. Robot Configuration and Operation Guide] -- [5.Enter Docker (For JETSON Nano and RPi 5)]**. For Orin boards, simply open a terminal and enter the commands mentioned in this course.

 This example uses `model: "qwen/qwen2.5-v1-72b-instruct:free", "qwen-v1-latest"`

 The responses from the large model may not be exactly the same for the same test command and may differ slightly from the screenshots in the tutorial. To increase or decrease the diversity of the large model's responses, refer to the section on configuring the decision-making large model parameters in the **[03.AI Model Basics] -- [5.Configure AI large model]** .

3. Running the Example

3.1 Starting the Program

For Raspberry Pi 5 and Jetson Nano controllers, you must first enter the Docker container. This is not necessary for the Orin board.

Open a terminal on the vehicle and enter the following command:

```
ros2 launch largemodel largemodel_control.launch.py text_chat_mode:=True
```

```
jetson@yahboom: ~  
-----  
jetson@yahboom:~$ ros2 launch largemodel largemodel control.launch.py text_chat_mode:=True  
[INFO] [launch]: All log files can be found below /home/jetson/.ros/log/2025-08-19-10-54-17-814305-yahboom-99292  
[INFO] [launch]: Default logging verbosity is set to INFO  
  
----- robot_type = A1, rplidar_type = c1, camera_type = usb -----  
  
-----robot_type = A1-----  
[INFO] [usb_cam_node_exe-1]: process started with pid [99356]  
[INFO] [joint_state_publisher-2]: process started with pid [99358]  
[INFO] [robot_state_publisher-3]: process started with pid [99360]  
[INFO] [Ackman_driver A1-4]: process started with pid [99362]  
[INFO] [base_node A1-5]: process started with pid [99364]  
[INFO] [imu_filter_madgwick node-6]: process started with pid [99366]  
[INFO] [ekf node-7]: process started with pid [99368]  
[INFO] [yahboom_joy A1-8]: process started with pid [99370]  
[INFO] [joy node-9]: process started with pid [99372]  
[INFO] [sllidar node-10]: process started with pid [99374]  
[INFO] [static_transform_publisher-11]: process started with pid [99376]  
[INFO] [model_service-12]: process started with pid [99393]  
[INFO] [action_service usb-13]: process started with pid [99409]  
[base_node A1-5] [INFO] [1755572058.715220575] [base_node]: Received parameters - linear_scale_x: 1.000000, linear_scale_y: 1.000000  
[static_transform_publisher-11] [WARN] [1755572058.731243357] []: Old-style arguments are deprecated; see --help for new-style arguments  
[imu_filter_madgwick node-6] [INFO] [1755572058.768449936] [imu_filter_madgwick]: Starting ImuFilter  
[sllidar node-10] [INFO] [1755572058.780476815] [sllidar_node]: SLLidar running on ROS2 package SLLidar.ROS2 SDK Version :1.0.1, SLLIDAR SDK Version:2.1.0  
[imu_filter_madgwick node-6] [INFO] [1755572058.783885747] [imu_filter_madgwick]: Using dt computed from message headers  
[imu_filter_madgwick node-6] [INFO] [1755572058.783980210] [imu_filter_madgwick]: The gravity vector is kept in the IMU message.  
[imu_filter_madgwick node-6] [INFO] [1755572058.786108865] [imu_filter_madgwick]: Imu filter gain set to 0.100000
```

Open another terminal and start the program:

```
ros2 run text_chat text_chat
```

```
jetson@yahboom: ~  
[System Information]  
IP Address 1: 192.168.11.198  
IP Address 2: 172.18.0.1  
-----  
ROS_DOMAIN_ID: 62 | ROS: humble  
my_robot_type: A1 | my_lidar: c1 | my_camera: usb  
-----  
jetson@yahboom:~$ ros2 run text_chat text_chat  
user input: 
```

3.2 Test Examples

Here are two reference test cases; users can create their own test commands.

- Please show me what items are in front of you and tell me what they are used for.
- Please look in front of you to see if there is a blue cube and a pack of tissues. If yes, nod your head. If not, shake your head.

⚠ Please do not include a period or any other characters at the end of the text you enter!

3.2.1 Example 1

Enter a test case in the text interface. After the model considers the answer, it responds to the user and performs actions according to the user's instructions.

```
jetson@yahboom:~$ ros2 run text_chat text_chat  
user input: Please show me what items are in front of you and tell me what they are used for.  
okay@, let me think for a moment... [INFO] [1755158911.300052374] [text_chat_node]: "action": ['seewhat()'], "response": Sure! I'm going to take a quick look around and see what's in front of me. Just a sec!  
user input: [INFO] [1755158918.419814638] [text_chat_node]: "action": [], "response": Alright, I've got a clear view of what's in front of me! There's a small fan on the desk, probably here to keep things cool during long work sessions. Next to it, I see a power strip with several chargers plugged in-perfect for keeping all your devices powered up. There's also a hand sanitizer bottle, which is great for staying clean and germ-free. A pen holder with pens and markers is nearby, so someone must be into writing or drawing. And of course, there are cables everywhere-typical office chaos! Let me know if you want me to focus on anything specific!  
[INFO] [1755158923.159358078] [text_chat_node]: "action": ['finishtask()'], "response": Got it! I've shared everything I could see and their uses. Let me know if you'd like to explore something else!  
End current task  
okay@, let me think for a moment... [INFO] [1755158935.086500530] [text_chat_node]: "action": ['finish_dialogue()'], "response": Alright, I'm stepping down for now. Have a wonderful day, and let me know if you need help later!  
user input: 
```

After completing a task, the robot enters a waiting state. Pressing the **ENTER** key allows you to continue the conversation. The command you enter is directly passed to the execution layer model, and all conversation history is retained. We can enter the "End Current Task" command again to end the robot's current task cycle and start a new one.

3.2.2 Example 2

Similar to Example 1, enter Example 2 in the terminal. The model will respond and execute the actions according to the command.

Example 2 uses the gimbal servo USB camera version as an example. Nodding and shaking the head can also be replaced with commands such as moving the robot forward or backward.

```
user input: Please look in front of you to see if there is a blue cube and a pack of tissues. If yes, nod your head. If not, shake your head.
okay0, let me think for a moment...-[INFO] [1755159997.734082084] [text_chat_node]: Decision making AI planning:1. Observe the current frame.
2. If both the blue cube and the pack of tissues are in the frame, call the function for the robotic arm to nod.
3. If either the blue cube or the pack of tissues is not in the frame, call the function for the robotic arm to shake its head.
user input: [INFO] [1755159000.891409449] [text_chat_node]: "action": ['seewhat()'], "arm_nod()'], "response": "I'm checking what's in front of me
right now, and if I see both the blue cube and the tissues, I'll give a little nod. If not, I'll shake my head to let you know!
[INFO] [1755159006.163648461] [text_chat_node]: "action": ['arm_shake()'], "response": "Hmm, I looked around carefully, but I couldn't spot the blu
e cube or the pack of tissues. Guess they're hiding somewhere else!
[INFO] [1755159010.678553180] [text_chat_node]: "action": ['arm_shake()'], "response": "I just checked again, and I still don't see the blue cube o
r the tissues. Maybe they're not here at all? I'm shaking my head to let you know!
[INFO] [1755159017.282908564] [text_chat_node]: "action": ['finishtask()'], "response": "Alright, I've checked and shaken my head-no blue cube or t
issues in sight. That's all for now!
[INFO] [1755159020.550603405] [text_chat_node]: "action": ['finishtask()'], "response": "I've already checked and confirmed-no blue cube or tissues
here. I'm all done for now!
End current task
okay0, let me think for a moment...-[INFO] [1755159059.809154571] [text_chat_node]: "action": ['finish_dialogue()'], "response": "Got it! I'm st
epping down and wrapping up this task. Have a great day!
user input: █
```

4. Source Code Analysis

Source code location:

Jetson Orin Nano:

```
#NUWA camera users
/home/jetson/yahboomcar_ros2_ws/yahboomcar_ws/src/largemodel/largemodel/action_s
ervice_nuwa.py
#USB camera users
/home/jetson/yahboomcar_ros2_ws/yahboomcar_ws/src/largemodel/largemodel/action_s
ervice_usb.py
```

```
/home/jetson/yahboomcar_ros2_ws/yahboomcar_ws/src/largemodel/largemodel/model_se
rvice.py
```

Jetson Nano, Raspberry Pi:

You need to first enter Docker.

```
#NUWA camera users
/root/yahboomcar_ros2_ws/yahboomcar_ws/src/largemodel/largemodel/action_service_
nuwa.py
#USB camera users
/root/yahboomcar_ros2_ws/yahboomcar_ws/src/largemodel/largemodel/action_service_
usb.py
```

```
/root/yahboomcar_ros2_ws/yahboomcar_ws/src/largemodel/largemodel/model_service.p
y
```

action_service.py

In the **init_ros_communication** method of the CustomActionServer class, a publisher for the topic "seewhat_handle" is created to publish visual processing signals.

```
# Create a publisher and publish the seewhat_handle topic
self.seewhat_handle_pub = self.create_publisher(String, 'seewhat_handle', 1)
```

- USB Camera, action_service_usb.py

Created a topic subscriber for /usb_cam/image_raw to receive the color image from the camera

```
#Image Topic Subscribers
self.subscription = self.create_subscription(
    Image, self.image_topic, self.image_callback, 1
)
```

- nuwa depth camera, action_service_nuwa.py

Created a topic subscriber for /ascamera_hp60c/camera_publisher/rgb0/image to receive the color image from the depth camera

```
#Image Topic Subscribers
self.subscription = self.create_subscription(
    Image, self.image_topic, self.image_callback, 1
)
```

The image_callback callback function continuously subscribes to the depth camera's color image.

When the action server calls the **seewhat** method, it publishes a signal to the seewhat_handle topic, notifying the model server's model_service node to upload an image to the multimodal large model.

```
def seewhat(self):
    self.save_single_image()
    msg = String(data="seewhat")
    self.seewhat_handle_pub.publish(
        msg
    ) # Normalize, publish the seewhat topic, and call the large model by
    model_service
```

This will call the **save_single_image** method to save a single image file, image.png.

Image save path:

Jetson Orin Nano, Jetson Orin NX:

```
/home/jetson/yahboomcar_ros2_ws/yahboomcar_ws/install/largemodel/share/largemodel/
resources_file/image.png
```

Jetson Nano and Raspberry Pi need to enter docker first:

```
/root/yahboomcar_ros2_ws/yahboomcar_ws/install/largemodel/share/largemodel/resou
rces_file/image.png
```

```

def image_callback(self, msg):#Image callback function
    self.image_msg = msg

def save_single_image(self):
    """
    保存一张图片 / Save a single image
    """
    if self.image_msg is None:
        self.get_logger().warning("No image received yet.") # No images
received yet...
        return
    try:
        # 将ROS图像消息转换为OpenCV图像 / Convert ROS image message to OpenCV image
        cv_image = self.bridge.imgmsg_to_cv2(self.image_msg, "bgr8")
        # 保存图片 / Save the image
        cv2.imwrite(self.image_save_path, cv_image)
        display_thread = threading.Thread(target=self.display_saved_image)
        display_thread.start()

    except Exception as e:
        self.get_logger().error(f"Error saving image: {e}") # Error saving
image...

```

model_service.py

In the **init_ros_communication** method of the LargeModelService class,
a subscriber to the topic "seewhat_handle" is created to receive visual processing signals.

```

#Create a Seewhat subscriber
self.seewhat_sub = self.create_subscription(String, 'seewhat_handle',
self.seewhat_callback,1)

```

When the seewhat_callback callback function receives the visual processing signal, it calls the **dual_large_model_mode** method, and then calls the **instruction_process** method to request inference and upload an image from the depth camera to the multimodal large model.

```

def seewhat_callback(self, msg):
    if msg.data == "seewhat":
        if (
            self.regional_setting == "China"
        ): # 在线模型推理方式: 决策层推理+执行层监督 / Online model inference method:
Decision layer reasoning + Execution layer supervision
            self.dual_large_model_mode(type="image")
        else:
            self.dual_large_model_international_model(type="image")

```

```

def instruction_process(self, prompt, type, conversation_id=None):
    """
    根据输入信息的类型（文字/图片），构建不同的请求体进行推理，并返回结果）
    Based on the type of input information (text/image), construct different
    request bodies for inference and return the result.
    """
    json_str = None
    prompt_seewhat = self.prompt_dict[self.language].get("prompt_2")

```

```

if self.regional_setting == "China":
    if type == "text":
        raw_content = self.model_client.multimodalinfer(prompt)
    elif type == "image":
        raw_content = self.model_client.multimodalinfer(
            prompt_seewhat, image_path=self.image_save_path
        )
    json_str = self.extract_json_content(raw_content)
elif self.regional_setting == "international":
    if type == "text":
        result = self.model_client.TaskExecution(
            input=prompt,
            map_mapping=self.map_mapping,
            language=self.language_dict[self.language],
            conversation_id=conversation_id,
        )
        if result[0]:
            json_str = self.extract_json_content(result[1])
            self.conversation_id = result[2]
        else:
            self.get_logger().info(f"ERROR:{result[1]}")
    elif type == "image":
        result = self.model_client.TaskExecution(
            input=prompt_seewhat,
            map_mapping=self.map_mapping,
            language=self.language_dict[self.language],
            image_path=self.image_save_path,
            conversation_id=conversation_id,
        )
        if result[0]:
            json_str = self.extract_json_content(result[1])
            self.conversation_id = result[2]
        else:
            self.get_logger().info(f"ERROR:{result[1]}")

```