

# **MicroViT and LLM Models**

## **Technical Documentation**

MicroVIT Robotics System

Generated: January 30, 2026

# Table of Contents

1. Overview
2. MicroViT Model (Vision Transformer)
3. LLM Model (Qwen2.5:0.5B)
4. Integration Architecture
5. Performance Characteristics
6. Configuration and Usage
7. Technical Specifications

# 1. Overview

The MicroVIT robotics system uses a two-stage AI pipeline for intelligent perception and natural language generation:

| Stage             | Model                     | Purpose                                         | Device      |
|-------------------|---------------------------|-------------------------------------------------|-------------|
| Vision Processing | MicroViT (MobileViT)      | Fast image understanding and feature extraction | Jetson Orin |
| Text Generation   | Qwen2.5:0.5B (via Ollama) | Natural language message generation             | Jetson Orin |

This hybrid approach combines the speed of lightweight vision models with the creativity of large language models, enabling real-time AI message generation on edge devices with limited computational resources.

## 2. MicroViT Model (Vision Transformer)

### 2.1 Introduction

MicroViT is a lightweight Vision Transformer architecture designed for edge devices. In this project, we use **MobileViT** (apple/mobilevit-small) as a proxy implementation since MicroViT is not available as a pre-trained model. MobileViT follows similar design principles: efficient attention mechanisms, hybrid CNN-ViT architecture, and optimization for mobile/edge deployment.

### 2.2 Architecture Details

| Component            | Specification                          |
|----------------------|----------------------------------------|
| Base Model           | MobileViTForImageClassification        |
| Model Name           | apple/mobilevit-small                  |
| Parameters           | ~5.6 million                           |
| Input Size           | 256x256 pixels                         |
| Feature Dimensions   | 320 (S1 variant)                       |
| Pre-training Dataset | ImageNet-1k (1000 classes)             |
| Architecture Type    | Hybrid CNN-ViT (MobileViT blocks)      |
| Attention Mechanism  | Efficient Single Head Attention (ESHA) |
| Model Size           | ~22 MB (FP32) / ~11 MB (FP16)          |

### 2.3 Variants

The implementation supports three MicroViT variants (S1, S2, S3) with different feature dimensions:

| Variant | Feature Dimensions | Channels        | Use Case                          |
|---------|--------------------|-----------------|-----------------------------------|
| S1      | 320                | [128, 256, 320] | Default - Balanced speed/accuracy |
| S2      | 448                | [128, 320, 448] | Higher accuracy, slower           |
| S3      | 512                | [192, 384, 512] | Highest accuracy, slowest         |

### 2.4 Processing Pipeline

#### Step 1: Image Preprocessing

- Decode base64 JPEG string to PIL Image
- Convert to RGB format
- Resize to 256x256 pixels (model input size)
- Normalize pixel values to [-1, 1] range
- Convert to PyTorch tensor [1, 3, 256, 256]

■■ Time: ~5-10ms

### **Step 2: Feature Extraction**

- Forward pass through MobileViT encoder
  - Extract features from last encoder layer (before classification head)
  - Global average pooling if needed
  - Extract CLS token or pooled features
- Time: ~9ms (GPU) / ~50-200ms (CPU)

### **Step 3: Classification**

- Forward pass through classification head
- Softmax to get class probabilities
- Extract top-5 predictions
- Map ImageNet class IDs to human-readable labels

### **Step 4: Feature-to-Text Conversion**

- Convert class probabilities to natural language
- Format: 'Primary object detected: {class} ({confidence}%)'
- Add secondary detections if confidence > 5%
- Add contextual descriptions based on detected class
- Output: Structured text description for LLM

## 2.5 ImageNet Class Recognition

MicroViT recognizes 1000 ImageNet classes. The system focuses on common objects relevant to robotics:

| Category               | Examples                                               |
|------------------------|--------------------------------------------------------|
| Traffic Infrastructure | stop_sign, traffic_light, fire_hydrant, parking_meter  |
| Vehicles               | bicycle, motorcycle, bus, train, truck, boat, airplane |
| People & Animals       | person, bird, cat, dog, horse, cow, elephant           |
| Furniture              | chair, couch, bed, dining_table, toilet                |
| Electronics            | tv, laptop, mouse, keyboard, cell_phone, remote        |
| Household Items        | bottle, cup, bowl, book, clock, vase, scissors         |
| Outdoor Objects        | bench, backpack, umbrella, suitcase, frisbee           |

## 2.6 Example Output

Input: Image of a turnstile/gate captured by Nano camera

**MicroViT Output:**

'Primary object detected: turnstile (45.2% confidence). Other possibilities: metal gate (12.3%), fence (8.1%), barrier (5.4%). Traffic infrastructure detected, road environment context.'

## 3. LLM Model (Qwen2.5:0.5B)

### 3.1 Introduction

Qwen2.5:0.5B is a lightweight Large Language Model developed by Alibaba Cloud. It is specifically designed for edge devices with limited computational resources. The model has approximately 0.5 billion parameters, making it suitable for real-time text generation on devices like Jetson Orin.

### 3.2 Model Specifications

| Property       | Value                          |
|----------------|--------------------------------|
| Model Name     | Qwen2.5:0.5B                   |
| Parameters     | ~0.5 billion (494.03M)         |
| Model Format   | GGUF (quantized)               |
| Quantization   | Q4_K_M (4-bit, medium quality) |
| Model Size     | ~380 MB (compressed)           |
| Context Length | 512 tokens (configurable)      |
| Max Tokens     | 500 tokens (configurable)      |
| Architecture   | Decoder-only Transformer       |
| Framework      | Ollama (local inference)       |
| Execution Mode | CPU-only (OLLAMA_NO_GPU=1)     |

### 3.3 Why Qwen2.5:0.5B?

#### Advantages for Edge Deployment:

- **Small Size:** ~380MB fits in limited RAM
- **Fast Inference:** ~200-500ms per generation on CPU
- **Good Quality:** Despite small size, generates coherent text
- **Multilingual:** Supports multiple languages
- **Quantized:** Q4\_K\_M quantization balances quality and speed
- **Ollama Integration:** Easy deployment via Ollama framework

### 3.4 Text Generation Process

#### Step 1: Prompt Construction

The system creates a structured prompt combining:

- MicroViT visual description
- LiDAR sensor data (distance, direction, confidence)
- Robot position and status
- Task instructions (generate status message)

## **Step 2: Ollama API Call**

- POST request to <http://localhost:11434/api/generate>
- Model: qwen2.5:0.5b
- Temperature: 0.8 (creative but focused)
- Max tokens: 500 (configurable)
- Stream: false (wait for complete response)

## **Step 3: Response Generation**

- Ollama processes prompt token by token
  - Generates natural language response
  - Returns complete message
- Time: ~3-5 seconds (CPU mode)

## **3.5 Example Prompt**

You are Robot robot1\_orin at location (0.00, 0.00). Based on my visual analysis (preprocessed with MicroViT), I can see: Primary object detected: turnstile (45.2% confidence). Other possibilities: metal gate (12.3%), fence (8.1%). Traffic infrastructure detected, road environment context. My LiDAR sensors report: - Distance to nearest obstacle: 0.5 meters - Obstacle direction: 0.0 degrees - Obstacle size: 0.5 meters - Sensor type: LiDAR (REAL) - Confidence level: 0.95 Generate a creative, informative status message (under 150 words) that: 1. Describes what I can see visually based on the preprocessed image features 2. Reports what my LiDAR sensors detect 3. Provides a combined assessment of the situation 4. Suggests what action I should take Make it sound like a robot reporting to other robots. Be engaging and specific about my location.

## 3.6 Example Generated Message

Hello everyone! Thank you for asking. Based on my visual analysis with MicroViT, I can see that one possible object in your environment is the turnstile. This is not an accident but rather an expected feature due to the distance from where I am to the nearest obstacle and the size of the obstacle. My LiDAR sensors detected an object within a 0.5-meter range with no obstacles around it. The sensor report indicates confidence at 95%, which is appropriate for detecting small objects like this. Based on my observations, I suggest that you take a closer look to see if there are any other interesting or unexpected features in the area where the turnstile might be located. If you encounter anything out of place or need additional assistance, please don't hesitate to let me know! Thank you for your understanding. Have a great day!

## 4. Integration Architecture

### 4.1 Two-Stage Pipeline

The system uses a two-stage approach to separate visual understanding from language generation:

| Stage | Component           | Input            | Output            | Time      |
|-------|---------------------|------------------|-------------------|-----------|
| 1     | Image Capture       | Camera feed      | Base64 JPEG       | ~10-20ms  |
| 2     | MicroViT Processing | Base64 image     | Text description  | ~50-200ms |
| 3     | Prompt Construction | MicroViT + LiDAR | Structured prompt | <1ms      |
| 4     | Ollama Generation   | Prompt           | Natural language  | ~3-5s     |
| 5     | MQTT Publishing     | AI message       | MQTT topic        | ~10-50ms  |

### 4.2 Data Flow

#### Nano → Orin (XML-RPC):

Camera image (base64 JPEG) → XML-RPC server → Orin AI service

#### Orin Internal Processing:

Base64 image → MicroViT → Feature extraction → Text description → Prompt construction  
→ Ollama → Natural language message

#### Orin → Controller (MQTT):

AI message → MQTT broker → Controller analysis

## 5. Performance Characteristics

### 5.1 Speed Comparison

| Model                | GPU Inference  | CPU Inference | Memory Usage | Accuracy        |
|----------------------|----------------|---------------|--------------|-----------------|
| MicroViT (MobileViT) | ~9ms           | ~50-200ms     | ~150MB       | Good (ImageNet) |
| BLIP (captioning)    | ~500ms         | ~2-5s         | ~1GB         | Excellent       |
| ViT-Base             | ~20ms          | ~300-500ms    | ~500MB       | Very Good       |
| Qwen2.5:0.5B         | N/A (CPU only) | ~200-500ms    | ~400MB       | Good            |

### 5.2 Total Pipeline Time

#### End-to-End Latency (CPU mode on Jetson Orin):

- Image capture: ~10-20ms
- MicroViT processing: ~50-200ms
- Ollama generation: ~3-5 seconds
- MQTT publishing: ~10-50ms

**Total: ~3.1-5.3 seconds per message**

With 30-second detection interval, this provides ample time for processing while maintaining real-time operation.

### 5.3 Memory Usage

#### Total System Memory (Jetson Orin):

- MicroViT model: ~150MB
- Qwen2.5:0.5B model: ~400MB
- System overhead: ~200MB
- Runtime buffers: ~100MB

**Total: ~850MB (well within 8GB RAM limit)**

# 6. Configuration and Usage

## 6.1 Environment Variables

| Variable              | Default               | Description                   |
|-----------------------|-----------------------|-------------------------------|
| USE_MICROVIT          | true                  | Enable/disable MicroViT       |
| MICROVIT_MODEL_NAME   | apple/mobilevit-small | Hugging Face model name       |
| MICROVIT_VARIANT      | S1                    | Variant: S1, S2, or S3        |
| MICROVIT_USE_CPU      | true                  | Force CPU mode                |
| OLLAMA_MODEL          | qwen2.5:0.5b          | Ollama model name             |
| OLLAMA_TEXT_MODEL     | qwen2.5:0.5b          | Text generation model         |
| OLLAMA_NO_GPU         | 1                     | CPU-only mode (avoid GPU OOM) |
| OLLAMA_MAX_TOKENS     | 500                   | Maximum tokens to generate    |
| OLLAMA_CONTEXT_LENGTH | 512                   | Context window size           |

## 6.2 Code Usage

```
# Initialize MicroViT from microvit_integration import MicroViTModel
model = MicroViTModel( model_name='apple/mobilevit-small', use_cpu=True, variant='S1' )
model.load_model() # Analyze image
image_description =
model.analyze_image(base64_image_string) # Returns: "Primary object detected:
turnstile (45.2% confidence)...". # Use with Ollama prompt = f"Based on my
visual analysis: {image_description}" ai_message = ollama.generate(prompt)
```

## 7. Technical Specifications

### 7.1 MicroViT Technical Details

| Specification            | Value                                  |
|--------------------------|----------------------------------------|
| Model Architecture       | MobileViT (Hybrid CNN-ViT)             |
| Attention Type           | Efficient Single Head Attention (ESHA) |
| Patch Size               | 16x16 pixels                           |
| Embedding Dimension      | 320 (S1)                               |
| Number of Layers         | Variable (MobileViT blocks)            |
| Activation Function      | GELU                                   |
| Normalization            | Layer Normalization                    |
| Optimizer (Training)     | AdamW                                  |
| Learning Rate (Training) | 1e-4                                   |
| Batch Size (Training)    | 256                                    |
| Precision                | FP32 (CPU) / FP16 (GPU)                |

### 7.2 Qwen2.5:0.5B Technical Details

| Specification     | Value                               |
|-------------------|-------------------------------------|
| Architecture      | Decoder-only Transformer            |
| Context Window    | 512 tokens (configurable up to 32K) |
| Vocabulary Size   | 151,936 tokens                      |
| Number of Layers  | 24                                  |
| Hidden Size       | 1,152                               |
| Attention Heads   | 12                                  |
| FFN Dimension     | 2,816                               |
| Activation        | SwiGLU                              |
| Position Encoding | RoPE (Rotary Position Embedding)    |
| Quantization      | Q4_K_M (4-bit, medium)              |
| Framework         | Ollama (GGUF format)                |

## Conclusion

The MicroViT robotics system successfully combines lightweight vision processing (MicroViT/MobileViT) with efficient language generation (Qwen2.5:0.5B) to enable real-time AI message generation on edge devices. This two-stage approach provides the optimal balance

between speed, accuracy, and resource constraints, making it ideal for autonomous robotics applications with limited computational resources.