



FPT University HCMC

Hỗ trợ robot hiểu ngôn ngữ con người và thực hiện các tác vụ đa dạng, phức tạp

Mentor: Anh Lương Hữu Dũng

Sinh viên thực hiện:

Đào Khang— SE183427

Huỳnh Khánh Quang— SE183777

—

MỤC LỤC

1. Lời nói đầu	3
2. Mục tiêu và Vấn đề	3
2.1. Mục tiêu	3
2.2. Vấn đề	3
3. Phương pháp	4
3.1. Model	4
3.1.1. YOLOv8x-seg	4
3.1.2. spacy	4
3.2. Danh sách các hàm chính	4
3.2.1. def speech_to_text()	4
3.2.2. def extract_object_and_location(user_input)	4
3.2.3. def detect_segment_object(image_path, target_object, target_location)	5
3.2.4. Main	5
4. Kết quả	6
4.1. Model gốc không finetune	6
4.2. Model đã finetune	7
4.2.1. Finetune bu lông	7
5. Kết luận và Hướng phát triển	9
5.1. Kết luận	9
5.2. Hướng phát triển	9
5.2.1. Về hướng phát triển:	9
5.2.2. Về các nội dung cần nghiên cứu:	9
5.2.3. Về các kết quả mong đợi:	9

DANH SÁCH HÌNH ẢNH

1. Kết quả của Model YOLOv8x-seg gốc dùng để nhận diện và phân đoạn vị trí trái cây (ở đây là chuối) theo vị trí trái, phải, giữa	6
2. Kết quả của Model YOLOv8x-seg gốc dùng để nhận diện và phân đoạn vị trí trái cây (ở đây là táo) theo vị trí trên, dưới	6
3. Dataset	7
4. Thông số trên từng epoch	7
5. Kết quả	8

1. Lời nói đầu

Trong bối cảnh cách mạng và AI phát triển nhanh chóng, robot không chỉ có khả năng thực thi nhiệm vụ vật lý, mà ngày càng được trang bị thêm khả năng nhận diện ngôn ngữ tự nhiên và tương tác linh hoạt với con người. Mục tiêu của dự án này là xây dựng một hệ thống robot có thể nhận lệnh ngôn ngữ tiếng Việt tự nhiên, dịch sang tiếng Anh, trích xuất thông tin về đối tượng và vị trí, sau đó nhận dạng và segment đối tượng trong ảnh dựa theo YOLOv8-seg để rồi đưa qua hệ thống Camera 3D và robot cụ thể ở đây là AUBO để thực hiện mệnh lệnh đã được đề ra

2. Mục tiêu và Vấn đề

2.1. Mục tiêu

- Xây dựng hệ thống nhận lệnh ngôn ngữ tiếng Việt tự nhiên.
- Dịch sang tiếng Anh để phù hợp với các mô hình xử lý NLP.
- Trích xuất đối tượng và vị trí mong muốn.
- Áp dụng YOLOv8 segmentation để nhận dạng và tô màu phân đoạn vật thể tương ứng trong ảnh và xác định trọng tâm vật ấy
- Kết hợp cùng với hệ thống Camera 3D và Robot để thực hiện nhiệm vụ

2.2. Vấn đề

- Tiếng Việt có ngữ pháp phức tạp đối với xử lý NLP.
- Trình phiên dịch Tiếng Anh còn nhiều cấu trúc chưa thể giống Tiếng Việt
- Kết hợp nhiều module: nhận diện giọng nói, dịch thuật, NLP, và Computer Vision.
- Dữ liệu ảnh thực tế có nhiều nhiễu tạp, vị trí đối tượng phức tạp.
- Chưa xử lý được những hình ảnh 3D

3. Phương pháp

3.1. Model

3.1.1. YOLOv8x-seg

Dự án sử dụng mô hình YOLOv8x-seg, là phiên bản lớn nhất của dòng mô hình YOLOv8 chuyên cho nhiệm vụ segmentation (phân vùng đối tượng trong ảnh). YOLOv8x-seg có khả năng phát hiện vật thể với độ chính xác cao, hỗ trợ segment chính xác từng pixel và hoạt động tốt trong thời gian thực.

3.1.2. spacy

Ngoài ra, dự án cũng sử dụng mô hình spaCy, là một mô hình NLP nhẹ nhưng hiệu quả, được huấn luyện để nhận biết cấu trúc ngữ pháp tiếng Anh như thực thể, cụm danh từ, động từ, giới từ, vị trí,...

3.2. Danh sách các hàm chính

3.2.1. def speech_to_text()

- **Mục đích:** Ghi âm giọng nói người dùng bằng micro, nhận diện nội dung tiếng Việt, sau đó dịch sang tiếng Anh.
- **Thư viện:** speech_recognition, deep_translator
 - **recognizer.listen():** Dùng để ghi âm
 - **recognizer.recognize_google():** Dùng để chuyển giọng nói thành văn bản tiếng Việt
 - **GoogleTranslator.translate():** Dùng để dịch văn bản tiếng Việt vừa nhận sang tiếng Anh

3.2.2. def extract_object_and_location(user_input)

- **Mục đích:** Dùng thư viện NLP spaCy để trích xuất thông tin từ câu văn: đối tượng (object) và vị trí (location).
- **Thư viện:** spacy
 - **nlp(user_input):** Dùng để phân tích cú pháp câu
 - **token.text.lower() in valid_locations:** Dùng để xác định từ chỉ vị trí
 - **doc.noun_chunks:** Dùng để tìm cụm danh từ để xác định object

3.2.3. def detect_segment_object(image_path, target_object, target_location)

- **Mục đích:** Nhận diện và phân đoạn đối tượng theo label và vị trí chỉ định trong ảnh bằng YOLOv8 segmentation.
- **Thư viện:** cv2, numpy, ultralytics, matplotlib
 - **YOLO(image):** phát hiện và phân đoạn đối tượng
 - **label == target_object:** Lọc ra đối tượng cần được phân đoạn
 - **cv2.fillPoly(), cv2.rectangle(), cv2.circle():** Dùng để hiển thị mask, bounding box, trọng tâm
 - **matplotlib.pyplot.imshow():** Hiển thị hình ảnh kết quả

3.2.4. Main

- **user_input = speech_to_text():** Ghi âm, dịch, và trả về câu tiếng Anh.
- **extract_object_and_location(user_input):** Trích object và location từ câu lệnh.
- **detect_segment_object(image_path, object_name, target_location):** Dựa vào object + location để nhận diện đúng vật thể.

4. Kết quả

4.1. Model gốc không finetune

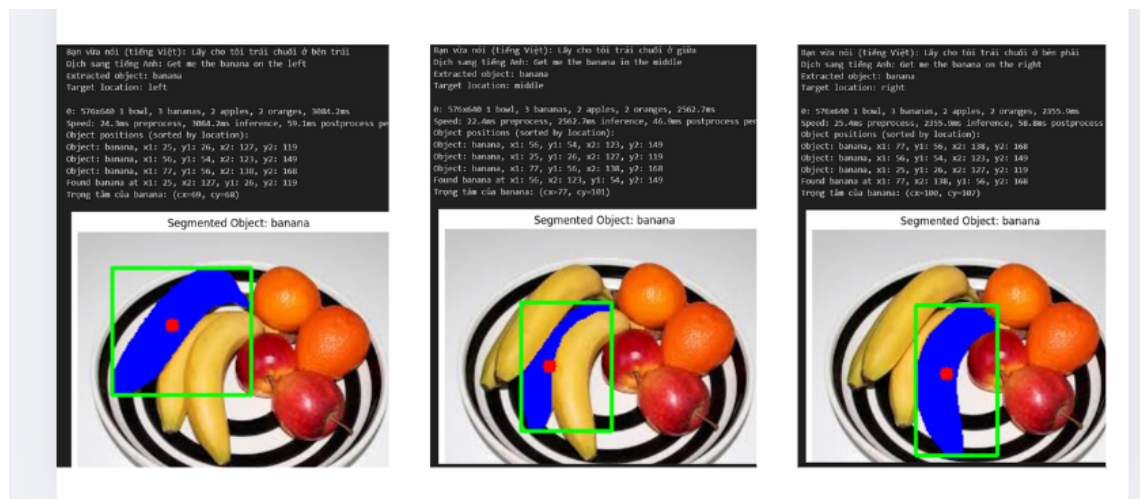


Figura 1: Kết quả của Model YOLOv8x-seg gốc dùng để nhận diện và phân đoạn vị trí trái cây (ở đây là chuối) theo vị trí trái, phải, giữa

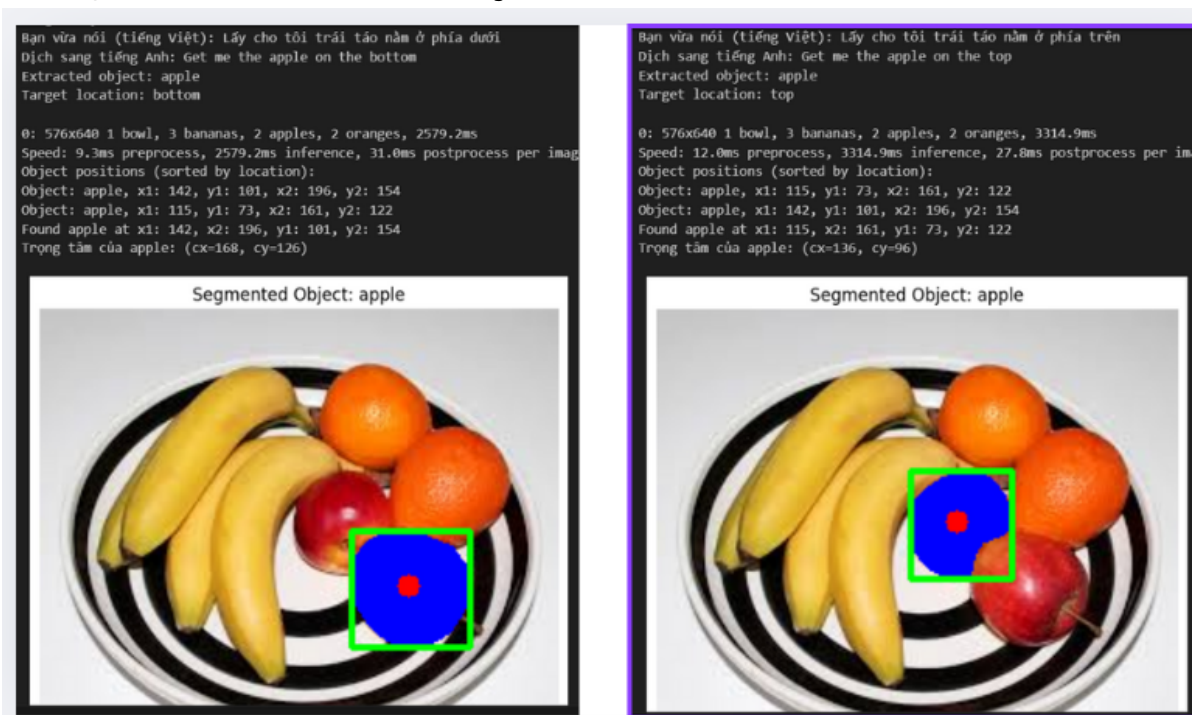


Figura 2: Kết quả của Model YOLOv8x-seg gốc dùng để nhận diện và phân đoạn vị trí trái cây (ở đây là táo) theo vị trí trên, dưới

4.2. Model đã finetune

4.2.1. Finetune bu lông

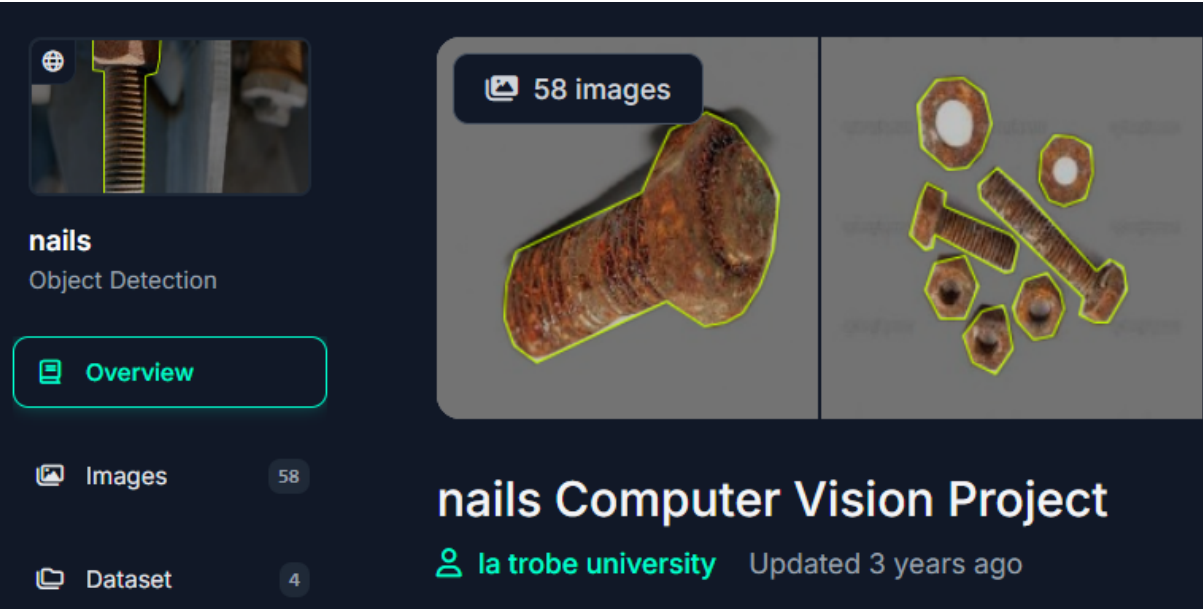


Figura 3: Dataset

Epoch	GPU_mem	box_loss	seg_loss	cls_loss	dfl_loss	Instances	Size	Mask(P)	10/10	[00:07:00:00, 1.37it/s]
1/50	6.85G	1.561	4.395	2.772	1.825	40	640: 100%	100%	10/10	[00:07:00:00, 1.37it/s]
	Class	Images	Instances	Box(P)	R	mAP50	mAP50-95)			1/1 [00:20:00:00, 20.89s/it]
Epoch	GPU_mem	box_loss	seg_loss	cls_loss	dfl_loss	Instances	Size	Mask(P)	10/10	[00:05:00:00, 1.74it/s]
2/50	10.4G	1.181	2.099	1.83	1.51	62	640: 100%	100%	10/10	[00:05:00:00, 1.74it/s]
	Class	Images	Instances	Box(P)	R	mAP50	mAP50-95)			1/1 [00:14:00:00, 14.93s/it]
Epoch	GPU_mem	box_loss	seg_loss	cls_loss	dfl_loss	Instances	Size	Mask(P)	10/10	[00:06:00:00, 1.65it/s]
3/50	10.5G	0.8966	1.587	1.275	1.289	43	640: 100%	100%	10/10	[00:06:00:00, 1.65it/s]
	Class	Images	Instances	Box(P)	R	mAP50	mAP50-95)			1/1 [00:16:00:00, 16.34s/it]

Epoch	GPU_mem	box_loss	seg_loss	cls_loss	dfl_loss	Instances	Size	Mask(P)	10/10	[00:05:00:00, 1.70it/s]
48/50	7.53G	0.3998	0.6235	0.2915	0.9382	14	640: 100%	100%	10/10	[00:05:00:00, 1.70it/s]
	Class	Images	Instances	Box(P)	R	mAP50	mAP50-95)			1/1 [00:00:00:00, 1.47it/s]
Epoch	GPU_mem	box_loss	seg_loss	cls_loss	dfl_loss	Instances	Size	Mask(P)	10/10	[00:05:00:00, 1.75it/s]
49/50	7.72G	0.3825	0.6065	0.2729	0.9103	21	640: 100%	100%	10/10	[00:05:00:00, 1.75it/s]
	Class	Images	Instances	Box(P)	R	mAP50	mAP50-95)			1/1 [00:00:00:00, 1.92it/s]
Epoch	GPU_mem	box_loss	seg_loss	cls_loss	dfl_loss	Instances	Size	Mask(P)	10/10	[00:06:00:00, 1.66it/s]
50/50	7.82G	0.3734	0.6227	0.2604	0.924	18	640: 100%	100%	10/10	[00:06:00:00, 1.66it/s]
	Class	Images	Instances	Box(P)	R	mAP50	mAP50-95)			1/1 [00:00:00:00, 1.97it/s]

Figura 4: Thông số trên từng epoch

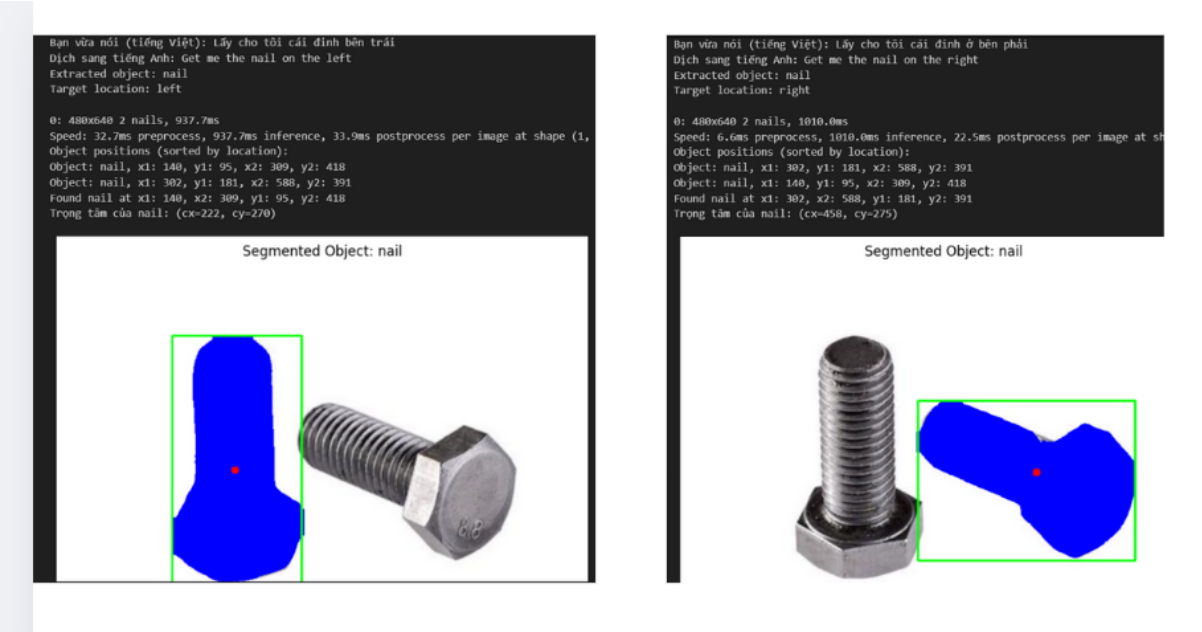


Figura 5: Kết quả

5. Kết luận và Hướng phát triển

5.1. Kết luận

Hệ thống đã có thể:

- Hiểu được câu lệnh bằng tiếng Việt.
- Trích xuất chính xác đối tượng (label) và vị trí (right, left, middle,...)
- Nhận dạng đúng đối tượng trong ảnh bằng YOLOv8-seg
- Phân đoạn đối tượng mong muốn

5.2. Hướng phát triển

5.2.1. Về hướng phát triển:

Về sự phát triển mà bọn em hướng tới, chúng em mong muốn tiếp tục cải tiến hệ thống bằng cách tích hợp một mô hình NLP mạnh mẽ hơn như các biến thể Transformer chuyên biệt để nâng cao khả năng hiểu ngữ nghĩa và xử lý câu lệnh phức tạp hơn. Đồng thời, chúng em cũng mong muốn được mở rộng dự án để có cơ hội được tiếp cận với hệ thống camera 3D và robot thực tế để có thể thử nghiệm mô hình trong môi trường thực, từ đó đưa giải pháp tiến gần hơn tới khả năng ứng dụng trong các bài toán robot gấp thà thông minh.

5.2.2. Về các nội dung cần nghiên cứu:

- Nghiên cứu về cách sử dụng và vận hành của Camera 3D
- Tìm hiểu thêm về lý thuyết cũng như cách kết nối và sử dụng Robot để có thể liên kết với Camera 3D và code Python

5.2.3. Về các kết quả mong đợi:

- Có thể kết nối hệ thống đã được lập trình trên Python với cái thiết bị tân tiến như Camera 3D để có thể xử lý hình ảnh, nhận diện và phân đoạn đối tượng cũng như kết nối với hệ thống Robot để có thể thực hiện những mệnh lệnh đã được đề ra bằng giọng nói trong thời gian thực