#### Integrated Vision-Physics-Reinforcement Learning Framework for Dynamic Industrial Robot Navigation







Benyamain yacoob author

yacoobby@udmercy.edu



xinyang zhang Author zhangxi24@udmercy.edu



Jingyuan wang Author wangji15@udmercy.edu



Mina maleki advisor malekimi@udmercy.edu

# Problem & Motivation

- Dynamic industrial environments pose challenges for mobile robot navigation due to unpredictable obstacles and changing layouts
- Traditional navigation methods struggle to adapt to these dynamic conditions without extensive reprogramming
- Recent research has explored adaptive techniques, integrating vision, physics, and reinforcement learning (RL)
- Our work addresses these challenges with a novel integrated framework combining vision-



## Related literature

- Recent research in robot navigation has focused on vision-based perception, physicsinformed control, and reinforcement learning approaches
- Vision-based navigation uses cameras and computer vision algorithms to perceive the environment and make navigation decisions
- Physics-informed control incorporates physical constraints and models into the control system to improve the robot's movement
- Reinforcement learning uses rewards and



#### Research Focus: CNN-PPO-Vision-Based Perception: CNN for object PINE And a stance estimation

- Integrated Framework: Combines CNN, PINN, and PPO for enhanced navigation in dynamic environments
- Physics-Informed Control: PINN for precise wheel dynamics control





#### Methodology: Integrated Vision, Physics, and RL Framework

- CNN: Detects target objects and estimates distances from RGB images and depth maps
- PINN: Controls wheel torque via physicallyconstrained dynamics
- PPO: Learns optimal navigation strategies through continuous state-action feedback



#### experimental environment



TurtleBot3 robotic platform used in our experiments





Gazebo: These components are demonstrated working in concert within a ROS2 Humble and Gazebo Classic simulation environment, showcasing the framework'



#### data Format



The dataset includes sensor readings LiDAR, RGB/Depth camera images. Using YOLOv516 model to generate bounding boxes and categories, and extended the label format to include target distances



3 0.322656 0.720833 0.107813 0.297222 4.69 3 0.897656 0.374306 0.087500 0.206944 10.00 1 0.392969 0.397917 0.093750 0.387500 7.55 3 0.254688 0.621528 0.037500 0.215278 5.89 2 0.127734 0.713889 0.253906 0.547222 4.50 1 0.666797 0.346528 0.030469 0.223611 10.00 1 0.809766 0.329861 0.067969 0.279167 10.00 4 0.916016 0.929861 0.166406 0.131944 3.58

<class\_id> <x\_center> <y\_center> <width> <heig

- All values except class\_id are normalized to [0, 1]
- Distance is in meters, representing camera-to-objection
- Format supports both object detection and distance

/nce>

enter

imation

# Problems with distance information processing



Direct depth estimation challenges showing inaccurate distance measurement (4.1 m) due to inclusion of multiple objects within the bounding box



#### Camera Top Angle Processing



Schematic diagram of the monocular camera ranging algorithm showing the pinhole imaging principle



# Cross-validation of distance information



Region growing segmentation results using uniform threshold, showing limitations in object boundary detection



# Cross-validation of distance information



Improved region growing segmentation using category-specific thresholds, demonstrating better object boundary definition



#### Distance information





Visualization of raw depth image data used for distance estimation



## Depth camera mapping attempt



Processing pipeline showing: (a) Original RGB image, (b) Segmented depth image, (c) Horizontal projection, and (d) Bounding box detection



# Cross-validation of distance information



Evaluating LiDAR recognition for target depth distance variation



#### dataset preprocessing

![](_page_15_Figure_1.jpeg)

![](_page_15_Picture_2.jpeg)

The training dataset comprises 10,077 images split in an 8:2 ratio between training and validation sets, with a batch size of 4 over 50 epochs

Preprocessed with Python scripts for format and class checks

![](_page_15_Picture_5.jpeg)

#### Vision-Based Perception (CNN)

- Modified Faster R-CNN with ResNet101 + FPN backbone
- Joint object detection & distance estimation
- Depth refinement using region-growing segmentation with class-specific thresholds for accurate distance measurement
- RGB + LiDAR data integrated for robust environmental understanding

![](_page_16_Picture_5.jpeg)

#### Faster R-CNN framework

![](_page_17_Figure_1.jpeg)

![](_page_17_Picture_2.jpeg)

## feature extraction

Faster R-CNN architecture (input  $3 \times 720 \times 1280$ )

Stages	Roles	Feature map changes (input $3 \times 720 \times 1280$ )
Conv1+bn+relu	Initial convolution + activation	[64, 360, 640]
Maxpool	Downsampling	[64, 180, 320]
Layer1	Residual module $\times$ 3	[256, 180, 320]
Layer2	Residual module $\times$ 4 + down-	[512, 90, 160]
	sampling	
Layer3	Residual module $\times$ 23 + down-	[1024, 45, 80]
	sampling	
Layer4	Residual module $\times$ 3 + down-	[2048, 23, 40]
	sampling	
Avgpool	Pooling to 1x1	[2048, 1, 1]
FC	Classification output	[1000]

$$\mathcal{L}_{ ext{total}} = \mathcal{L}_{ ext{cls}} + \mathcal{L}_{ ext{bbox}} + lpha \cdot \mathcal{L}_{ ext{dist}}$$

$$\mathcal{L}_{ ext{dist}} = rac{1}{N} \sum_{i=1}^{N} \left| d_i^{ ext{pred}} - d_i^{ ext{gt}} 
ight|$$

Our model is based on Faster R-CNN with a ResNet-FPN backbone. The standard classification and bounding box heads are retained, using crossentropy and Smooth L1 loss respectively. We additionally introduce a custom distance regression head that outputs the estimated object-camera

Stage 2

![](_page_18_Picture_6.jpeg)

#### Physics-Informed Control (PINN)

Neural network embeds wheel dynamics:

- where J represents the wheel's moment of inertia, b is the damping coefficient, θ is the wheel angle, τ is the motor torque and the load torque
- Inputs: Wheel states, torques, timestep (8D input vector)

![](_page_19_Picture_4.jpeg)

#### Physics-Informed Control (PINN)

- [64, 32] hidden layers with ReLU activations
- Multi-threaded ROS2 services for real-time torque updates at 30 Hz
- Replaces speed-control plugins with reverse torque strategy for direct control of wheel acceleration

![](_page_20_Picture_4.jpeg)

## **PINN** framework

![](_page_21_Figure_1.jpeg)

PINN torque control flowchart showing the multi-threaded control architecture

![](_page_21_Picture_3.jpeg)

### Reinforcement Learning (PPO)

 Agent Input (6D state): Target distance, obstacle distance, linear & angular velocity, visibility flag, normalized timestep

 $R_t = k_p (d_{t-1} - d_t) + k_v v_t - k_\omega \omega_t^2 + k_s \max(0, s_{\min,t} - s_{th}) + \mathbf{1}_{\text{goal}} R_{\text{goal}} - \mathbf{1}_{\text{coll}} R_{\text{coll}} + \mathbf{1}_{\text{goal}} R_{\text{goal}} - \mathbf{1}_{\text{coll}} R_{\text{goal}} - \mathbf{1}_{\text{goal}} R_{\text{goal}} - \mathbf{1}_{\text{coll}} R_{\text{goal}} - \mathbf{1}_{\text{coll}} R_{\text{goal}} - \mathbf{1}_{\text{coll}} R_{\text{goal}} - \mathbf{1}_{\text{goal}} R_{\text{go$ 

 Balances goal-seeking, obstacle avoidance, smooth motion, and energy usage

![](_page_22_Picture_4.jpeg)

### **Results & discussion**

- exhibited effective navigation in dynamic environments, with the PPO algorithm successfully learning to combine information from the CNN and PINN components to generate appropriate navigation commands
- Performance evaluation showed consistently high success rates in reaching target locations while maintaining smooth trajectory generation
- The system demonstrated reliable obstacle

![](_page_23_Picture_4.jpeg)

#### Optimizer Comparison

![](_page_24_Figure_1.jpeg)

![](_page_24_Figure_2.jpeg)

SGD

#### Adam

![](_page_24_Figure_5.jpeg)

#### SGD with Momentum

	AdamW	SGD	SGD + Momentum
Avg Loss	0.0612	0.5181	0.1442
Precision	0.8711	0.7536	0.8689
Recall	0.9348	0.9054	0.9369
Remarks	Fast convergence, strong on regression	Slow, weak performance	Balanced performance, strong generalization

AdamW converges the fastest with the lowest loss and highest precision among all optimizers.

23

Do

SGD converges slowly and shows unstable learning, with the lowest precision and highest loss.

SGD with Momentum achieves

### Yolov5l6 vs Faster RCNN (YOLOV5l6)

![](_page_25_Figure_1.jpeg)

- One-stage model, very fast
- Bounding boxes less accurate, sometimes overlapping
- No distance estimation
- Good for real-time tasks, but less reliable in complex scenes

![](_page_25_Picture_6.jpeg)

### Yolov5l6 vs Faster RCNN (faster rcnn)

![](_page_26_Figure_1.jpeg)

- Two-stage model, more accurate
- Boxes are tighter and more precise
- Supports real-world distance output
- Better at detecting occluded or cluttered objects

![](_page_26_Picture_6.jpeg)

# Ppo solution

![](_page_27_Figure_1.jpeg)

- Policy collapse due to unstable updates
- Delay or noise in visual perception (e.g., Model/Camera latency)
- Sparse or overly complex reward design

![](_page_27_Picture_5.jpeg)

# Ppo solution

![](_page_28_Figure_1.jpeg)

- Improve reward shaping and feedback clarity
- Integrate more stable perception models
- Fine-tune PPO hyperparameters (learning rate, entropy, batch size)

![](_page_28_Picture_5.jpeg)

### PINN result

![](_page_29_Figure_1.jpeg)

![](_page_29_Picture_2.jpeg)

#### Conclusion

- a modified Faster R-CNN architecture with ResNet101 backbone for effective object detection and distance estimation
- a PINN-based wheel dynamics controller operating at 30 Hz with precise torque control
- a PPO implementation with verified hyperparameters managing real-time navigation decisions

![](_page_30_Picture_4.jpeg)

#### Future work

- focus on bridging the sim-to-real gap through improved physics modeling
- optimizing computational efficiency for realtime operation, expanding the framework's capabilities to support a broader range of environments
- incorporating additional sensor modalities for enhanced perception

![](_page_31_Picture_4.jpeg)

![](_page_32_Picture_0.jpeg)

## questions?