# Secure Autonomy for Contested Environments

## Miroslav Pajic

CPSL@Duke

Department of Electrical and Computer Engineering

Department of Computer Science

Department of Mechanical Engineering and Material Science

Duke University

Duke PRATT SCHOOL of ENGINEERING

DUKE ROBOTICS

# Offline Reinforcement Learning with Off-Policy Evaluation



- Off-policy evaluation (OPE) is important for **filling the gap between training offline reinforcement learning (RL) controllers and choosing which one to deploy online**

| Platform Testing | Platform Not Available | Platform Not Available | Platform Not Available | Next Platform Test |
|---|---|---|---|---|

**Phase *I -- Testing***
Data collection using the latest controller.

**Phase *II – Offline RL***
Use the updated experience dataset to fine-tune existing controllers or train new ones.

**Phase *III -- OPE***
Estimate the controller candidates *offline* and select the one with best performance.
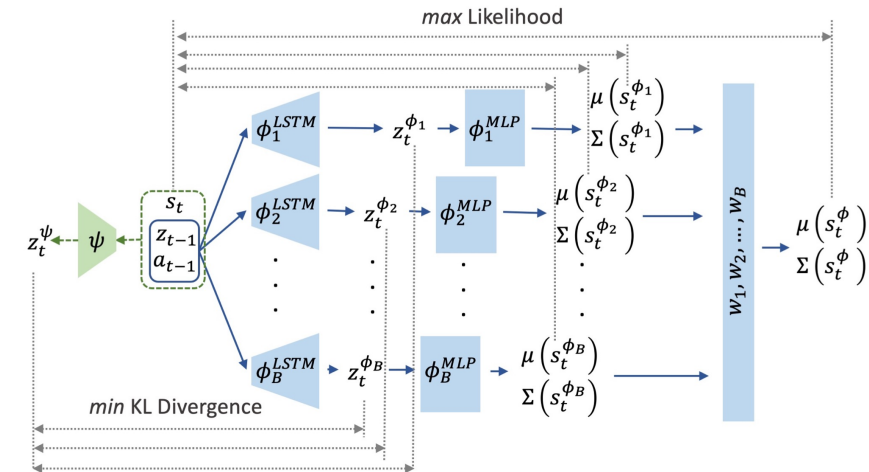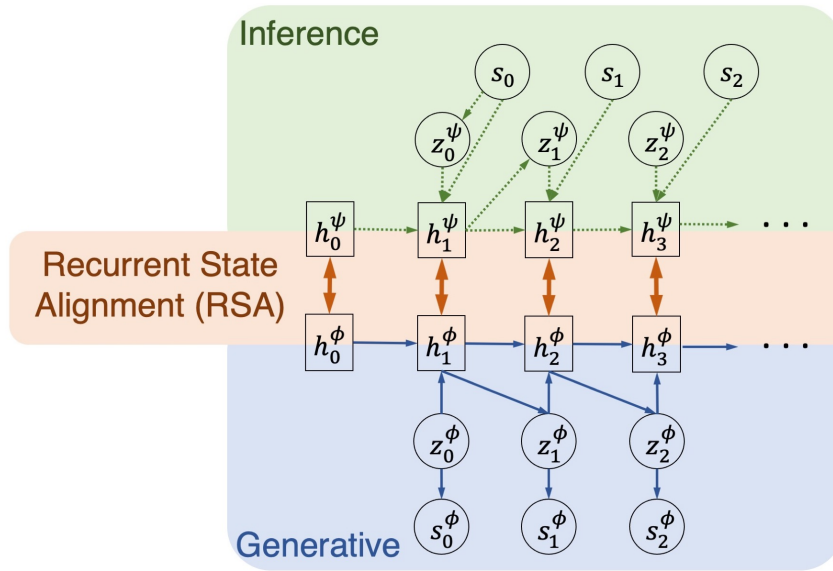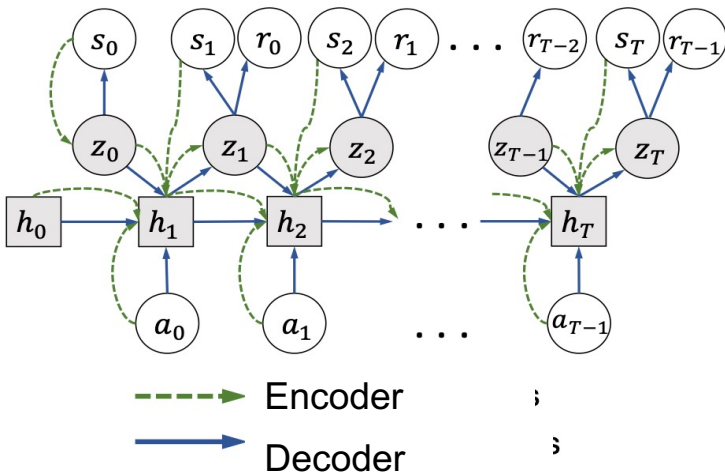
**Phase *IV -- Testing***
Run the best performing controller selected, and collect new trajectories.

# Variational Latent Branching Model (VLBM) for OPE (ICLR23)

- Formulate a latent space where latent variables can transit over time $p_\phi(z_t|z_{t-1}, a_{t-1})$

- Both encoder and decoder are **LSTMs**

- The encoder infuse the knowledge of the environment into the latent space

- The decoder generates synthetic trajectories **over time**

- Recurrent state alignment (RSA)
  - To mitigate the effect that decoder starts working **long after the encoder encodes the entire trajectory**
  - Minimize the **mean pairwise error** between LSTM states of encoder and decoder

- Branching for the decoder
  - Multiple decoders sample from the encoder to reduce variabilities possibly caused by, e.g., **random** initialization and **stochasticity** during training

- Overall training objective (maximize)
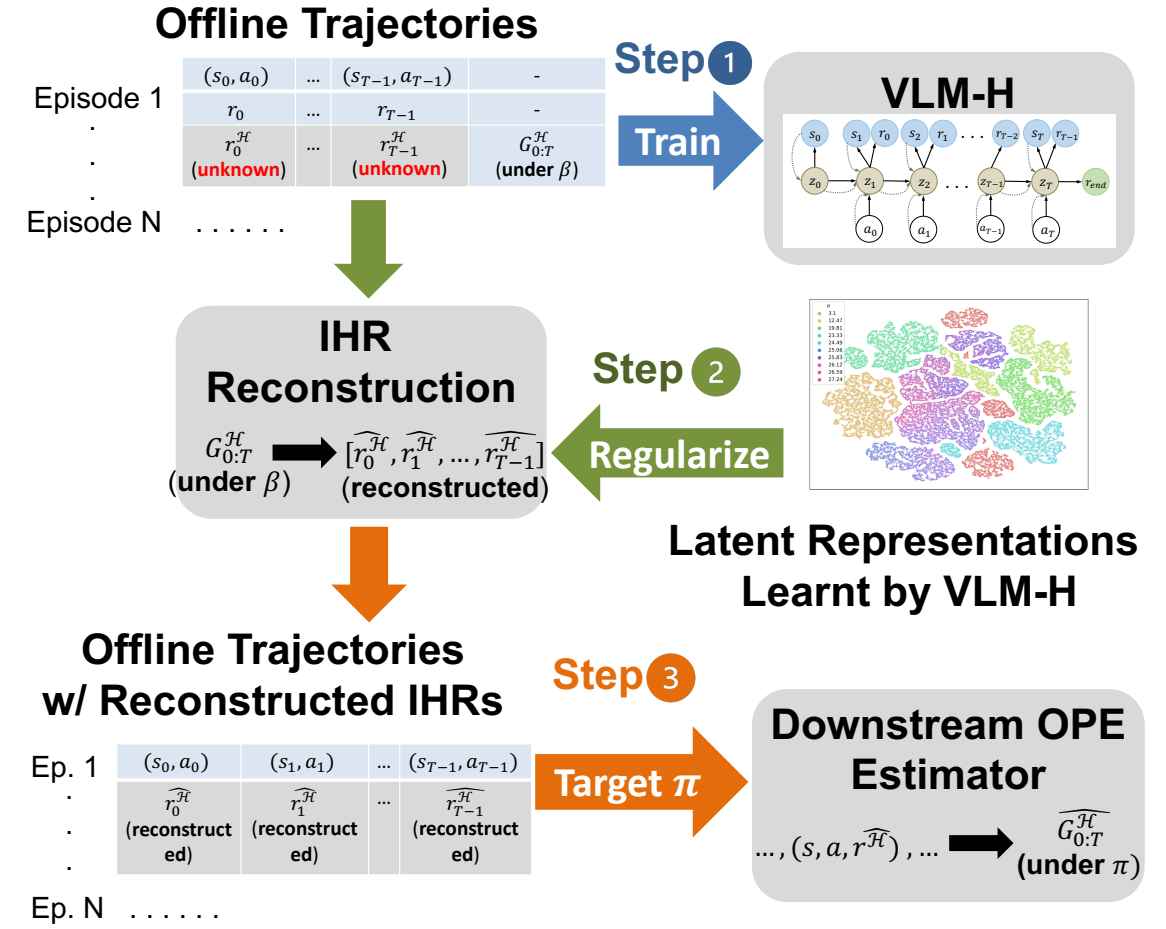
  *ELBO − RSA +*
  *log_likelihood_for_each_branch*

# Off-Policy Evaluation for Sparse/Human Feedback (NeurIPS23)

Unknown **Immediate Human**

...at the IHRs, $r_t^{\mathcal{H}}$,

...vable. Instead, the

...uman return,

...$_t^{\mathcal{H}}$, is available **at**
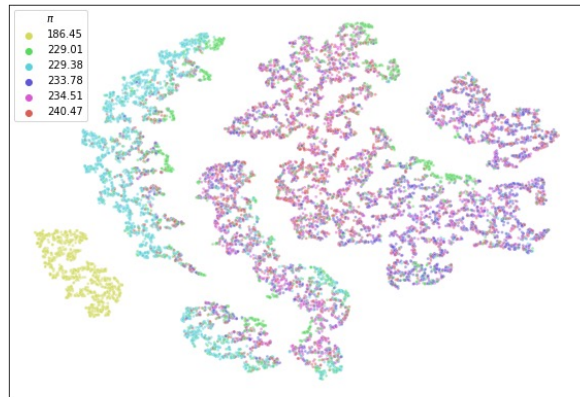
**...ch episode (i.e.,**

**...arse)**.

**Objective:** Given **a fixed set of offline trajectories collected by a behavioral policy** $\beta$, estimate the **expected total human return** over the unknown state-action visitation distribution $\rho^\pi$ of the target (evaluation) policy $\pi$ -- $\mathbb{E}_{(s,a)\sim\rho^\pi}[\sum_t \gamma^t r_t^{\mathcal{H}}]$.
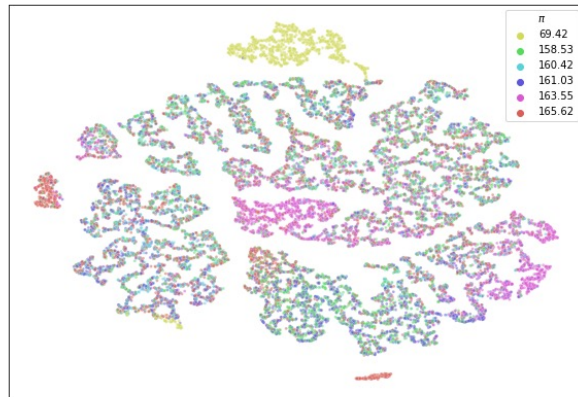
Duke
PRATT SCHOOL *of*
ENGINEERING

Unknown **Immediate Human**

...nat the IHRs, $r_t^{\mathcal{H}}$,
...vable. Instead, the

**Offline Trajectories**

Episode 1

| $(s_0, a_0)$ | ... | $(s_{T-1}, a_{T-1})$ | - |
|---|---|---|---|
| $r_0$ | ... | $r_{T-1}$ | - |
| $r_0^{\mathcal{H}}$ (unknown) | ... | $r_{T-1}^{\mathcal{H}}$ (unknown) | $G_{0:T}^{\mathcal{H}}$ (under $\beta$) |

Episode N    . . . . . .

**Step** 1

**Train**

**VLM-H**

$s_T$  $r_{T-}$

$z_T$

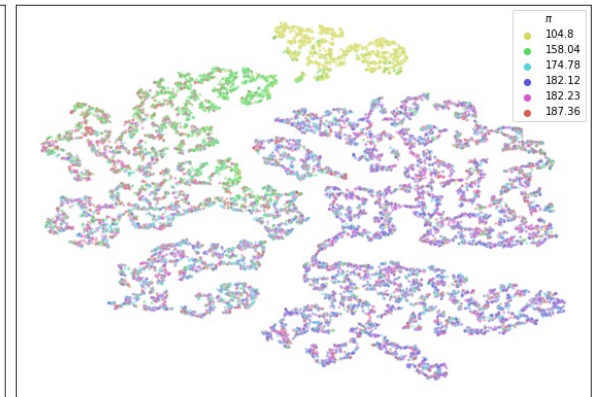$a_T$

coding)

ecoding]



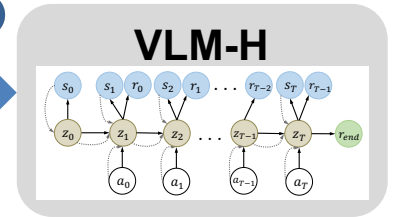Patient #0            Patient #1            Patient #2            Patient #3

**trajectories collected by a behavioral policy** $\beta$, estimate the **expected total human return** over the unknown state-action visitation distribution $\rho^\pi$ of the target (evaluation) policy $\pi$ -- $\mathbb{E}_{(s,a)\sim\rho^\pi}[\sum_t \gamma^t r_t^{\mathcal{H}}]$.

Ep. 1

| $(s_0, a_0)$ | $(s_1, a_1)$ | ... | $(s_{T-1}, a_{T-1})$ |
|---|---|---|---|
| $r_0^{\widehat{\mathcal{H}}}$ (reconstructed) | $r_1^{\widehat{\mathcal{H}}}$ (reconstructed) | ... | $r_{T-1}^{\widehat{\mathcal{H}}}$ (reconstructed) |

Ep. N   . . . . . .

**Target** $\pi$

**Estimator**

$..., (s, a, r^{\widehat{\mathcal{H}}}), ... \Longrightarrow \widehat{G_{0:T}^{\mathcal{H}}}$ (under $\pi$)

Unknown **Immediate Human**

...nat the IHRs, $r_t^{\mathcal{H}}$, ...vable. Instead, the ...uman return, ...$_t^{\mathcal{H}}$, is available **at** ...ch episode (i.e., ...arse)**.

$s_T$  $r_{T-1}$

$z_T$  →  $r_{end}$

$a_T$

...coding)

...ecoding)

**Objective:** Given **a fixed set of offline trajectories collected by a behavioral policy** $\beta$, estimate the **expected total human return** over the unknown state-action visitation distribution $\rho^\pi$ of the target (evaluation) policy $\pi$ -- $\mathbb{E}_{(s,a)\sim\rho^\pi}[\sum_t \gamma^t r_t^{\mathcal{H}}]$.

**Offline Trajectories**

**Step ①**

Episode 1

| $(s_0, a_0)$ | ... | $(s_{T-1}, a_{T-1})$ | - |
| $r_0$ | ... | $r_{T-1}$ | - |
| $r_0^{\mathcal{H}}$ (unknown) | ... | $r_{T-1}^{\mathcal{H}}$ (unknown) | $G_{0:T}^{\mathcal{H}}$ (under $\beta$) |

Episode N    . . . . . .

**Train**

**VLM-H**

**Latent Representations Learnt by VLM-H**

**Step ②**

**IHR Reconstruction**

$G_{0:T}^{\mathcal{H}}$ (under $\beta$)  ⟶  $[\widehat{r_0^{\mathcal{H}}}, \widehat{r_1^{\mathcal{H}}}, ..., \widehat{r_{T-1}^{\mathcal{H}}}]$ (reconstructed)

**Regularize**

**Offline Trajectories w/ Reconstructed IHRs**

**Step ③**

Ep. 1

| $(s_0, a_0)$ | $(s_1, a_1)$ | ... | $(s_{T-1}, a_{T-1})$ |
| $\widehat{r_0^{\mathcal{H}}}$ (reconstructed) | $\widehat{r_1^{\mathcal{H}}}$ (reconstructed) | ... | $\widehat{r_{T-1}^{\mathcal{H}}}$ (reconstructed) |

Ep. N   . . . . . .

**Target $\pi$**

**Downstream OPE Estimator**

$..., (s, a, \widehat{r^{\mathcal{H}}}), ...$  ⟶  $\widehat{G_{0:T}^{\mathcal{H}}}$ (under $\pi$)
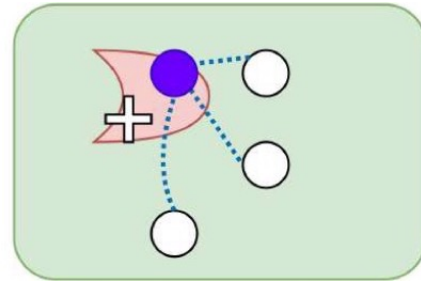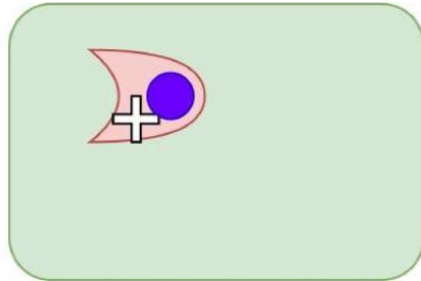
**Max-Min Optimization**

$$\max_{\theta \in \Theta} \min_{\phi \in \boxed{\Phi}} R(\theta, \phi)$$



- Inner minimization problem is difficult to solve → local-optimum
- Worst-case optimization can be over-conservative for *unrealistic* adversary (i.e., overly capable)

**Max-Min Optimization**

$$\max_{\theta \in \Theta} \min_{\phi \in \hat{\Phi}} R(\theta, \phi) = \max_{\theta \in \Theta} \min_{\phi_1, ..., \phi_m \in \hat{\Phi}} \min_{\phi \in \{\phi_i\}_{i=1}^m} R(\theta, \phi)$$
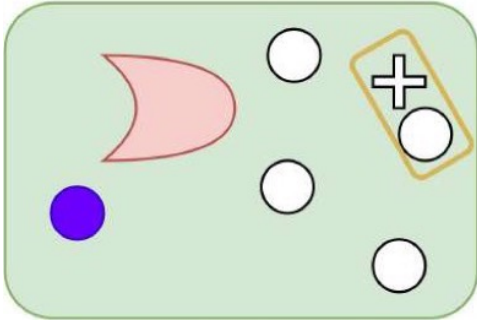


Learners
instead of
fixed
adversaries

***Efficient approximation of the inner optimization*** i.e., the size of adversary herd is upper-bounded to obtain sufficient approximation precision.

**Max-Min Optimization with Adversarial Herd – Optimization Over Worst-*k* Adversaries**

$$\max_{\theta \in \Theta} \min_{\phi \in \Phi} R(\theta, \phi) \quad \Longrightarrow \quad \max_{\theta \in \Theta} \min_{\phi_1, \ldots, \phi_m \in \hat{\Phi}} \frac{1}{I_{\theta, \hat{\Phi}, k}} R(\theta, \phi_i)$$
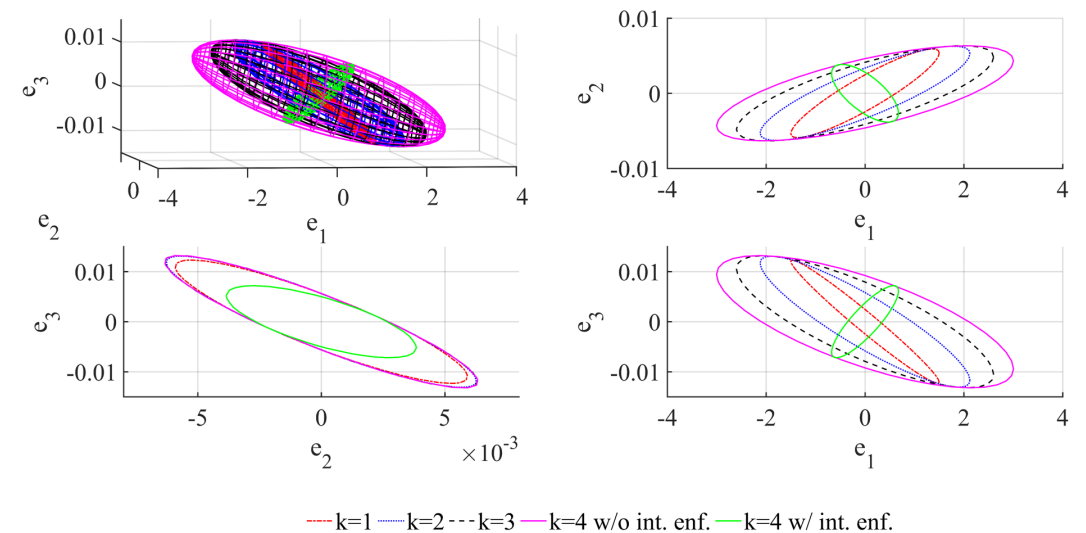


***Resolving Potential Over-Pessimism***
i.e., modify the objective from optimizing its worst-case performance,
to optimizing its average performance over the worst-k adversaries

If we choose a set of adversaries that are different enough, then the number of adversaries needed to approximate the inner optimization problem is in linear order of the desired precision.

If our objective is to use adversarial herd to approximate accurately with high probability, instead of an almost sure approximation, then the number of required adversaries can be reduced.

## Max-Min Optimization with Adversarial Herd – Optimization Over Worst-*k* Adversaries

$$\max_{\theta \in \Theta} \min_{\phi \in \Phi} R(\theta, \phi) \implies \max_{\theta \in \Theta} \min_{\phi_1, ..., \phi_m \in \hat{\Phi}} \frac{1}{I_{\theta, \hat{\Phi}, k}} R(\theta, \phi_i)$$



(a) Hopper  (b) Walker2d  (c) Half-Cheetah



(a) Baseline (0 adv)  (b) RARL (1 adv)  (c) RAP (population)  (d) ROLAH

(e) Baseline (0 adv)  (f) RARL (1 adv)  (g) RAP (population)  (h) ROLAH

(i) Baseline (0 adv)  (j) RARL (1 adv)  (k) RAP (population)  (l) ROLAH

# How can we analyze the impact of different attack vectors on CPS (i.e., QoC)?

An attack sequence

- is **strictly stealthy** iff

$$KL\big(Q(Y_{-\infty}^{-1}, Y_0^a : Y_t^a) || P(Y_{-\infty} : Y_t)\big) = 0$$

for any $t \geq 0$,

- is $\epsilon$-**stealthy** if

$$KL\big(Q(Y_{-\infty}^{-1}, Y_0^a : Y_t^a) || P(Y_{-\infty} : Y_t)\big) \leq \log(\frac{1}{1-\epsilon^2}) \text{ for any } t \geq 0.$$



The system is $(\epsilon, \alpha)$-attackable for arbitrarily large $\alpha$ and arbitrarily small $\epsilon$, if the closed-loop dynamics is incrementally exponentially stable (IES) in the set $S$ and the open loop dynamics is incrementally unstable in the set $S$.
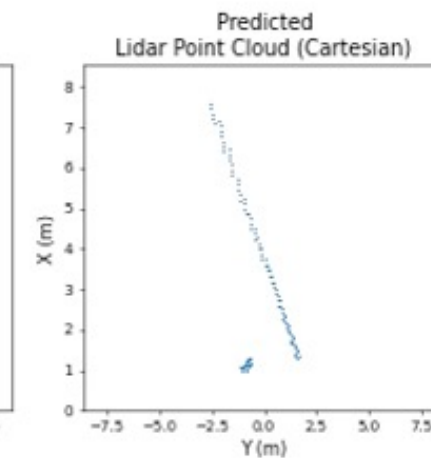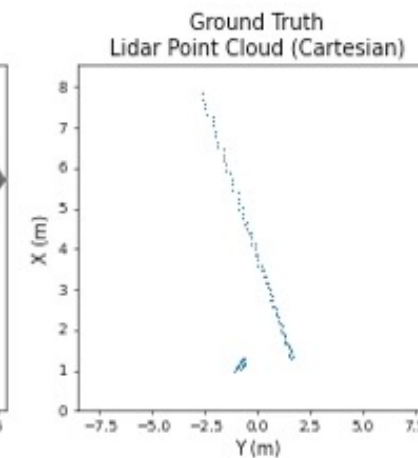
**Translation Attacks**



Attack Timeline



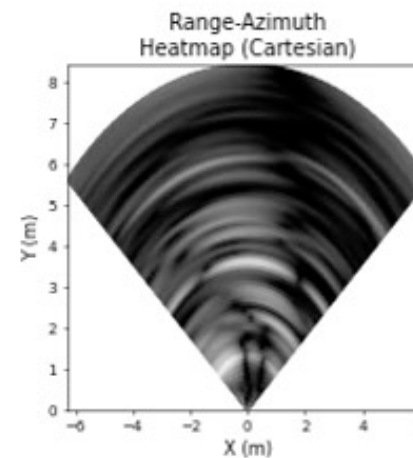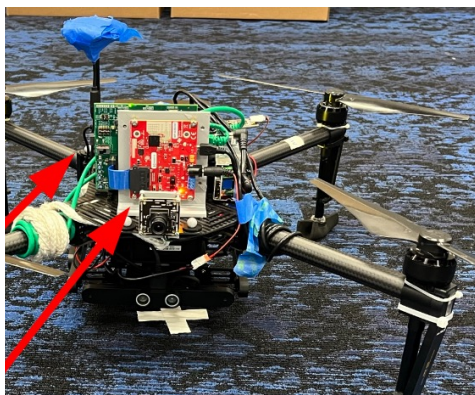https://sites.google.com/view/madradar-public/home
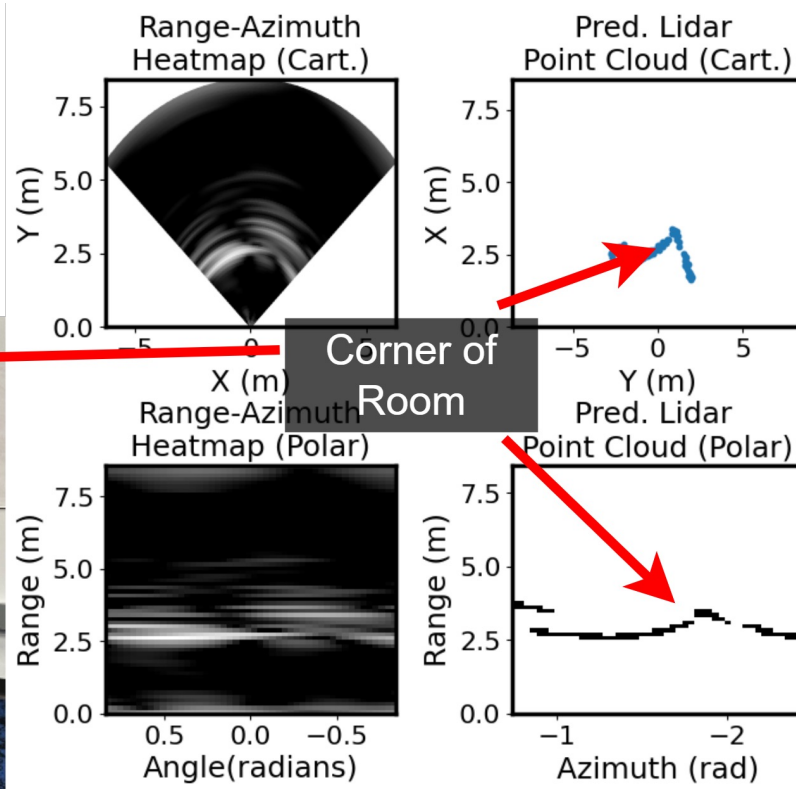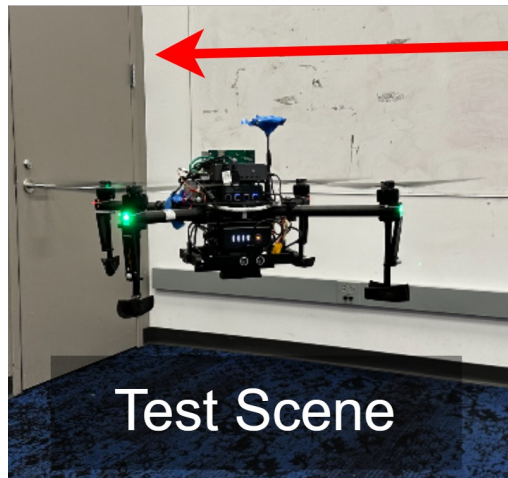
# mmWave-based Autonomy (ICRA'24*)

**Goal:** ***Low-cost*** (~$100), ***low-weight*** solution for
*adversarially robust* situational awareness and
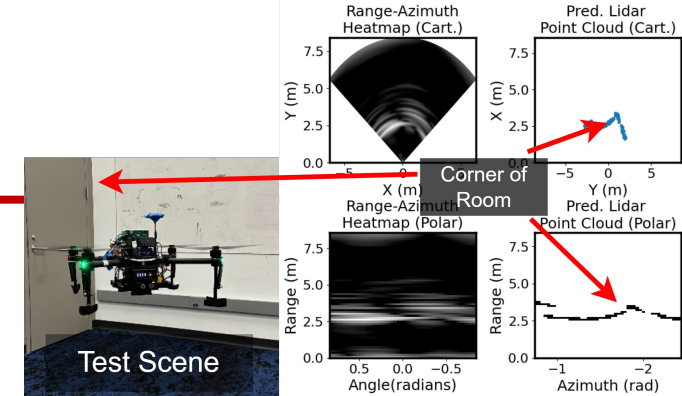autonomy on ***computationally constrained*** devices

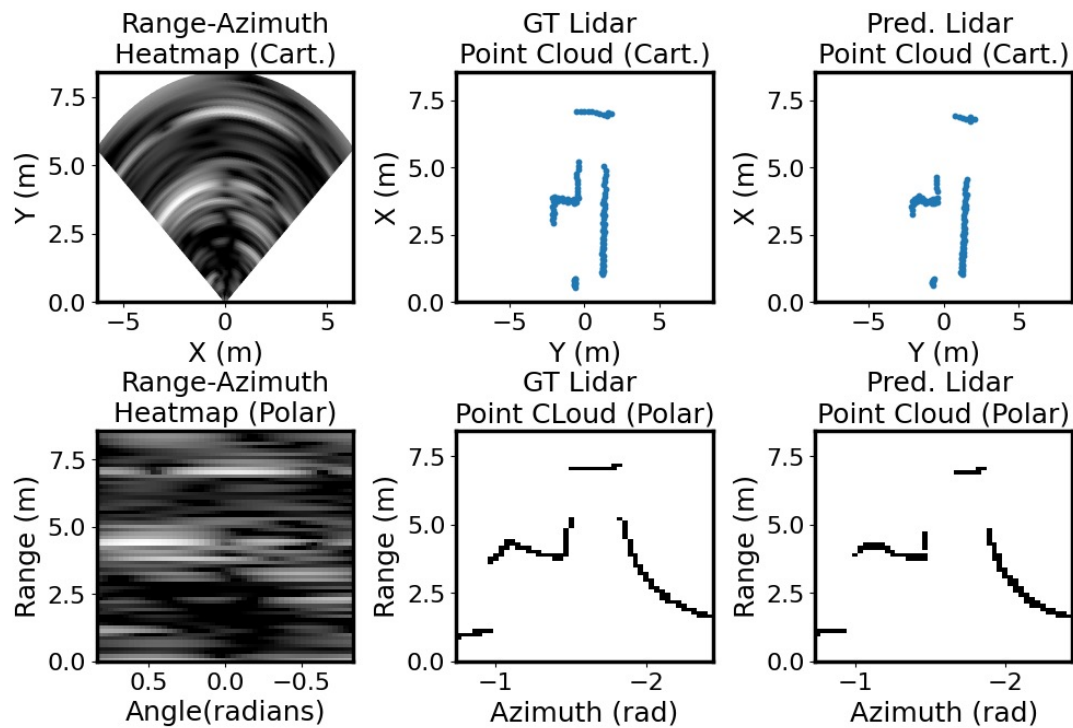- ***Real-time*** *high-accuracy* point clouds on a NUC-powered
  drone using mmWave sensing



Lobbies

Hallways

Laboratories

Office Spaces

Range-Azimuth Heatmap (Cartesian)

Ground Truth Lidar Point Cloud (Cartesian)

Predicted Lidar Point Cloud (Cartesian)

Ground Truth Lidar Point CLoud (Spherical)

Predicted Lidar Point Cloud (Spherical)

Nominal Operation

Rapid Movements

RadCloud Significantly Better

cham. RadCloud
mhaus. RadCloud
cham. 20 frames
mhaus. 20 frames
cham. 40 frames
mhaus. 40 frames

Test Scene

Range-Azimuth Heatmap (Cart.)

Pred. Lidar Point Cloud (Cart.)

Corner of Room

Range-Azimuth Heatmap (Polar)

Pred. Lidar Point Cloud (Polar)

# mmWave-based Autonomy (ICRA'24*)



Test Scene

Corner of Room

Range-Azimuth Heatmap (Cart.)

Pred. Lidar Point Cloud (Cart.)

Range-Azimuth Heatmap (Polar)

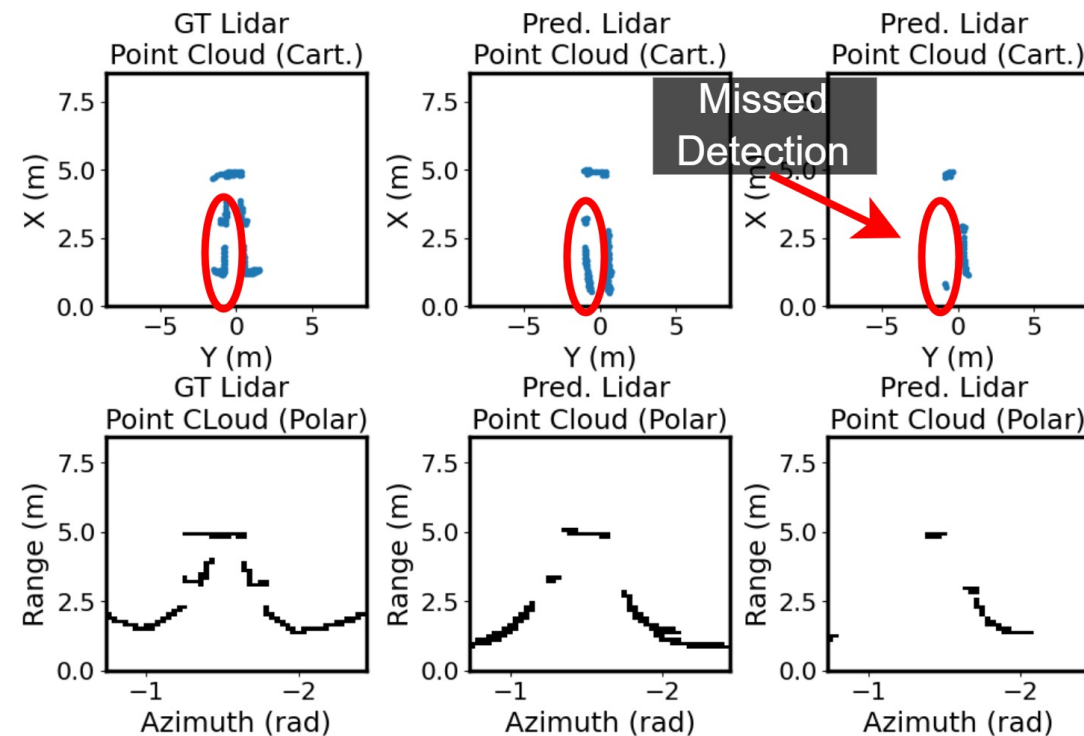Pred. Lidar Point Cloud (Polar)

## New Environments



Input radar data, ground truth point cloud, and predicted point cloud
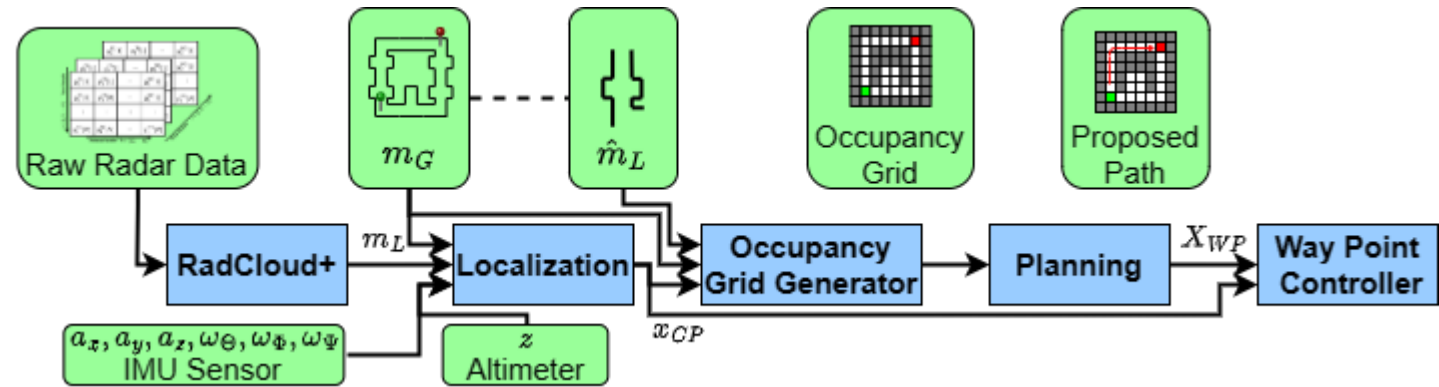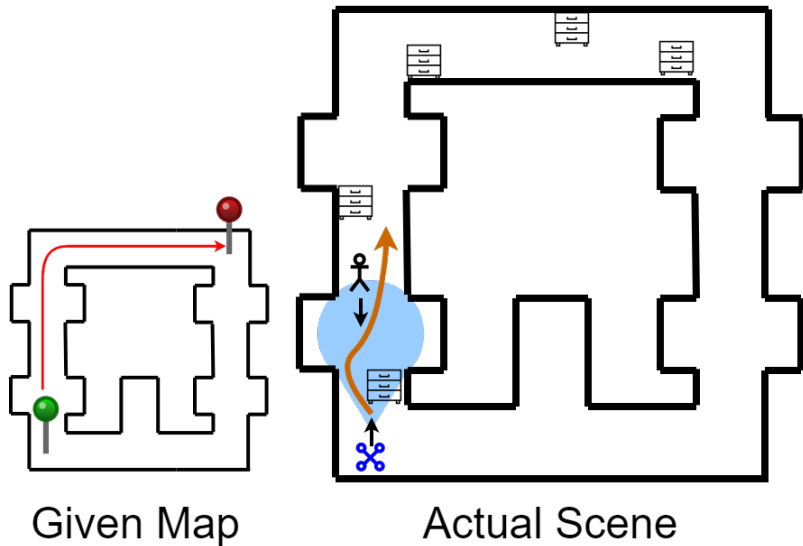
## Aggressive Maneuvers



Missed Detection

GT Lidar          RadCloud          40 Frames

# RadNav

- *Goal: N*avigate through an environment using only radar and traditional odometry sensors



Given Map     Actual Scene

# Thank you