# Structural Alignment in Worst-case Security Analysis and Multi-agent Design
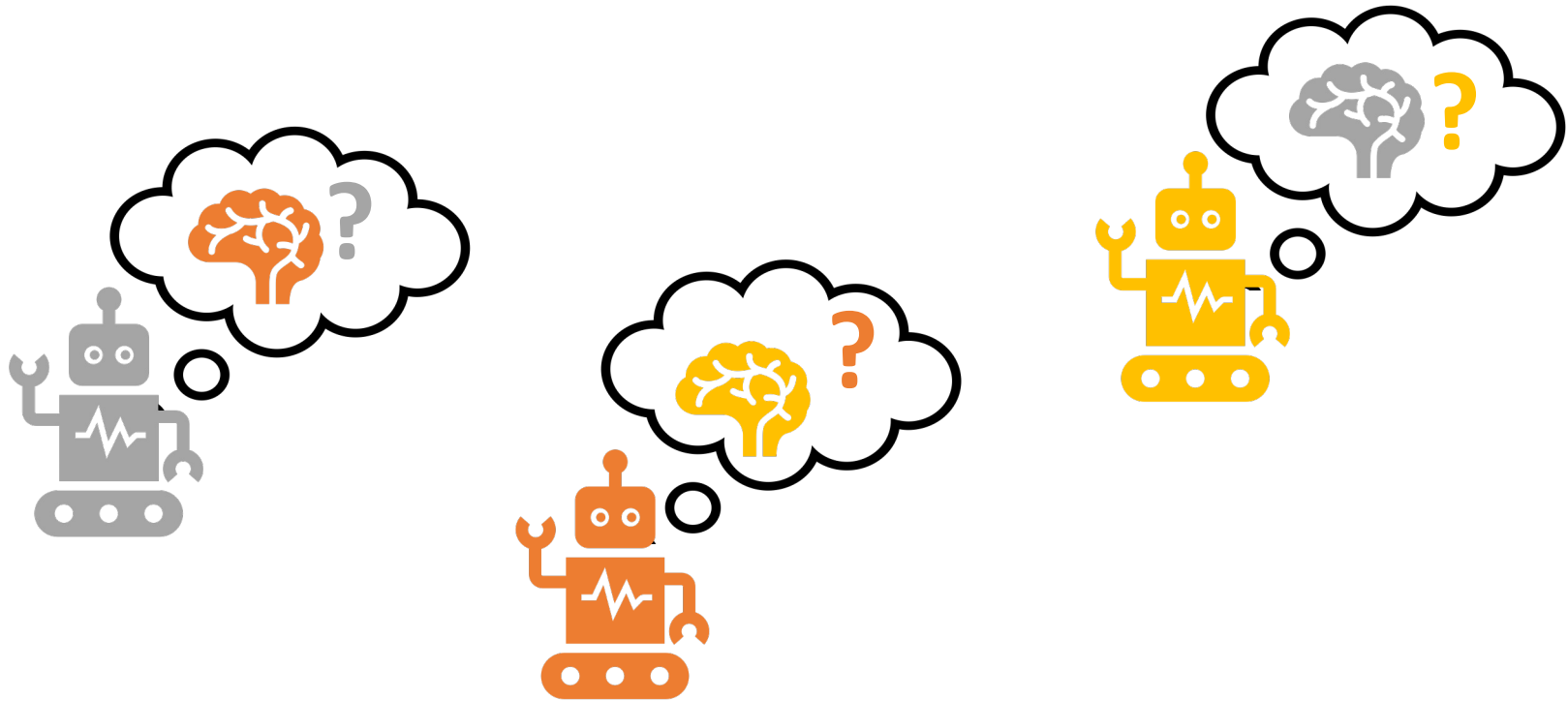
**Washington Garcia (UF)**

Kevin Butler (UF)

Pin-Yu Chen (IBM)

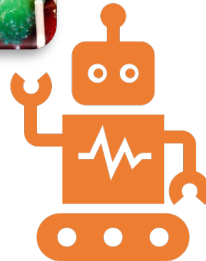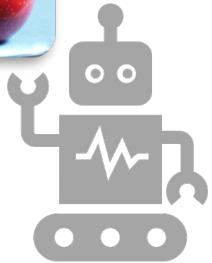Somesh Jha (UW)

Scott Clouse (AFRL/ACT3)

Entities in a multi-agent system must be aware of their surroundings, but also the *representational structure* of their counterparts. This can bolster situational awareness:
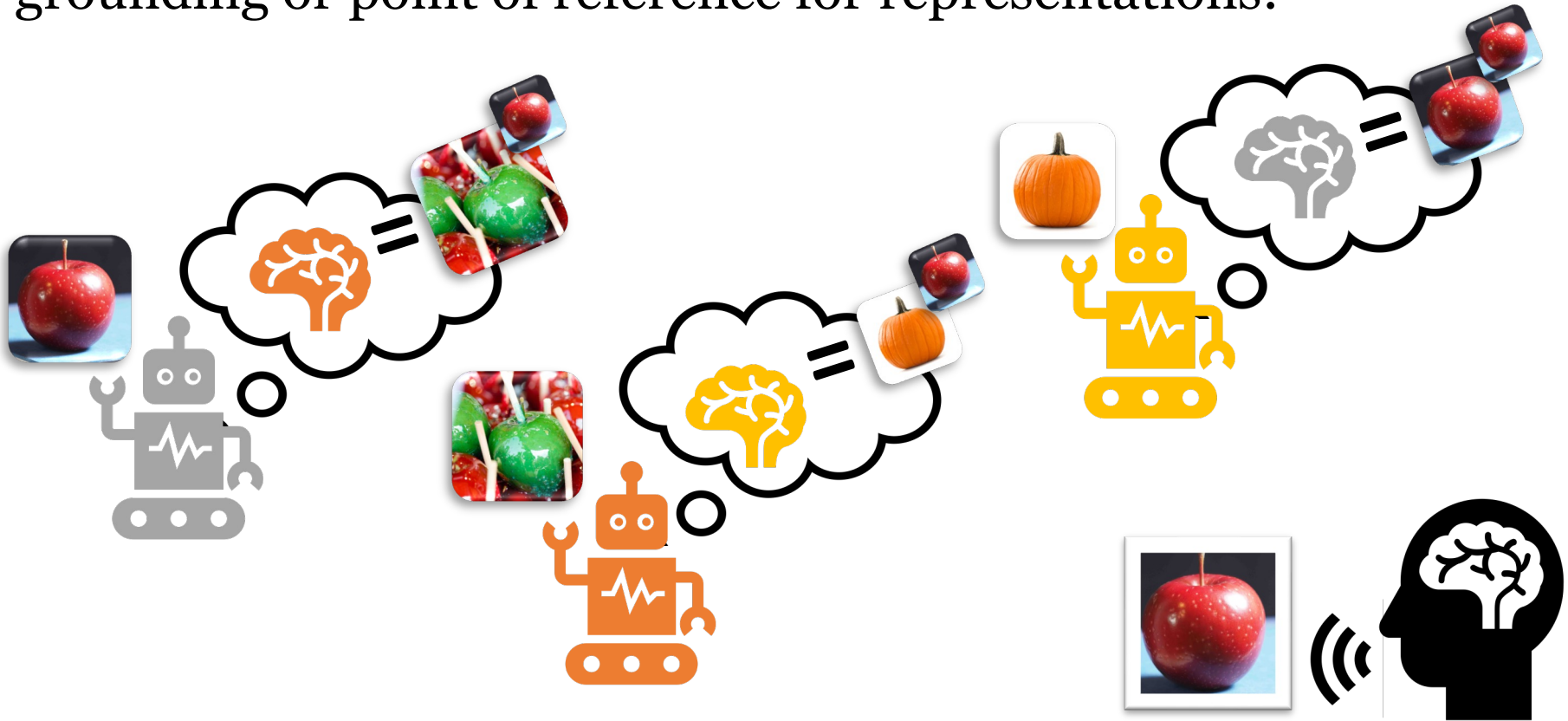
What is a good yardstick for representations?

A "human-in-the-loop" or proxy thereof can provide a grounding or point of reference for representations:

What is a good yardstick for representations?

A "human-in-the-loop" or proxy thereof can provide a grounding or point of reference for representations:

Two main questions:

1. Is agent concept resolution possible from an ***adversary's*** perspective? What's the worst-case security analysis?

2. It isn't always realistic to have a human-in-the-loop, or fine-grained data labels. What are ways around this?

A1: Expansion of our previous hard-label paper, in submission to IEEE SaTML 2023.

A2: Presented initial idea @ NAACL 2022. Introduced expansion during previous meeting. Then worked on it as part of AFRL 2022 internship.

Two main questions:

1. Is agent concept resolution possible from an ***adversary's*** perspective? What's the worst-case security analysis?

2. It isn't always realistic to have a human-in-the-loop, or fine-grained data labels. What are ways around this?

A1: Expansion of our previous hard-label paper, in submission to IEEE SaTML 2023.

A2: We presented an initial idea @ NAACL 2022. Introduced expansion during previous meeting. Then worked on it as part of AFRL 2022 internship.
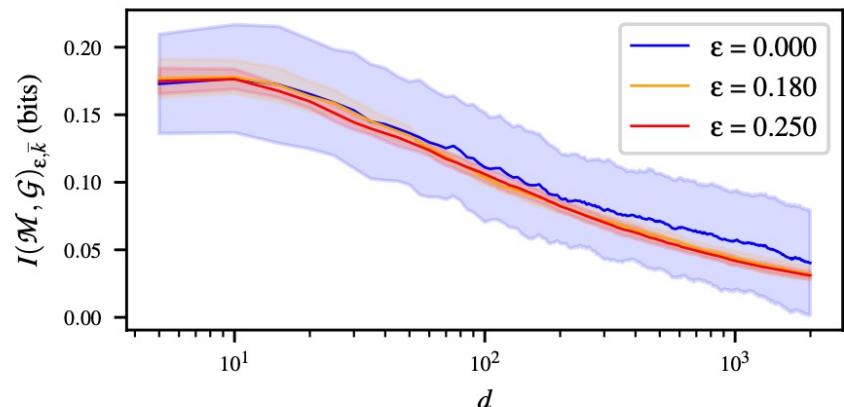
Assumptions:

1. Agent data lives on a low-dimensional manifold

2. The distribution of data points (M), true gradients (G), and ZO gradient estimates form a Markov chain:

$$\mathcal{M} \rightarrow \mathcal{G} \rightarrow \tilde{\mathcal{G}}$$

We showed previously that manifold-gradient *mutual information* can be modeled as a function of data dimension:

$$I(\mathcal{M}; \mathcal{G})_{\epsilon,k} = 2 \int_{\mathcal{M}^+} p(1, \mathbf{x}^+) \log(\frac{p(1, \mathbf{x}^+)}{p_{\mathcal{G}}(1) p_{\mathcal{M}}(\mathbf{x}^+)}) \, d\mathbf{x}^+$$
$$+ 2 \int_{\mathcal{M}^+} p(-1, \mathbf{x}^+) \log(\frac{p(-1, \mathbf{x}^+)}{p_{\mathcal{G}}(-1) p_{\mathcal{M}}(\mathbf{x}^+)}) \, d\mathbf{x}^+.$$

In practice, *does the Markov chain (MC) exist?*

First, show that the MC can be modeled by a zeroth-order (i.e., hard-label) adversary through two algorithms.

Local step neighborhood analysis (Algorithm 1):

**Algorithm 1:** Local Markov chain step (MC_step)

**Input:** Hard-label Gaussian process (GP), LIME kernel width $k$

**Output:** Sample feature coefficients $W \in \mathbb{R}^d$ and their quality score $R^2 \in \mathbb{R}$, GP result (res)

1 initialize LIME Ridge regression trainer (LIME) [31]
2 /* Execute GP to collect samples */
3 $X, Y, \text{res} \leftarrow \text{GP}()$

4 $f_W \leftarrow \text{LIME}(X, \ddot{Y}, k)$
5 $R^2 \leftarrow f_W(X)$
6 **return** $W, R^2, \text{res}$

Build local neighborhood from ZO queries as Gaussian process (GP) →

Use queries to train linear model of decision boundary. →

In practice, *does the Markov chain (MC) exist?*

First, show that the MC can be modeled by a zeroth-order (i.e., hard-label) adversary through two algorithms.

Local step neighborhood analysis (Algorithm 1):

Build local neighborhood from ZO queries as Gaussian process (GP)

Use queries to train linear model of decision boundary.



**Algorithm 1:** Local Markov chain step (MC_step)

Using Algorithm 1 to model local ZO queries, we can model the whole-attack agent Markov chain through Bayesian optimization (OptiLIME).

**Algorithm 2:** Markov chain probing of hard-label attack

1   $GP := (init, \mathbf{x}_0)$
2   kernel width $k \leftarrow$ OptiLIME(MC_step, GP)
3   /* Initialize through MC_step */
Call Algorithm 1  →  4   $W_{init}, R^2, \mathbf{x} \leftarrow$ MC_step(GP, $k$)
5   **for** $i := 1$ *to* $n$ **do** *Hard-label attack loop*
6     $GP \leftarrow$ (approximate_gradient, $\mathbf{x}$)
7     $k \leftarrow$ OptiLIME(MC_step, GP, $\mathbf{x}'$)
8     /* Approximate through MC_step */
Call Algorithm 1  →  9     $W_{\hat{\mathbf{g}}_i}, R^2, \boldsymbol{\theta} \leftarrow$ MC_step(GP, $k$)
10    Update $\mathbf{x}$ from $\boldsymbol{\theta}$ using attack formulation
11   **end**
Get "quality" of MC models ($R^2$)  →  12   **return** $\mathbf{x}$, $\{W_{init}, W_{\hat{\mathbf{g}}_1}, \ldots, W_{\hat{\mathbf{g}}_n}\}$,
     $\{R^2_{init}, R^2_{\hat{\mathbf{g}}_1}, \ldots, R^2_{\hat{\mathbf{g}}_n}\}$

Using Algorithm 1 to model local ZO queries, we can model the whole-attack agent Markov chain through Bayesian optimization (OptiLIME).



Call Algorithm 1 ⟶

Call Algorithm 1 ⟶

Get "quality" of MC models ($R^2$) ⟶

**Algorithm 2:** Markov chain probing of hard-label attack

1   GP := (init, $\mathbf{x}_0$)
2
3   $\bigcirc\!\!\bigcirc \rightarrow f_W \rightarrow R^2_{init}$
4
5   **for** $i := 1$ *to* $n$ **do** *Hard-label attack loop*
6     GP ← (approximate_gradient, $\mathbf{x}$)
7
8     $\bigcirc\!\!\bigcirc \rightarrow f_W \rightarrow R^2_{\hat{\mathbf{g}}_i}$
9
10     Update $\mathbf{x}$ from $\theta$ using attack formulation
11   **end**
12   **return** $\mathbf{x}$, $\{W_{init}, W_{\hat{\mathbf{g}}_1}, \ldots, W_{\hat{\mathbf{g}}_n}\}$,
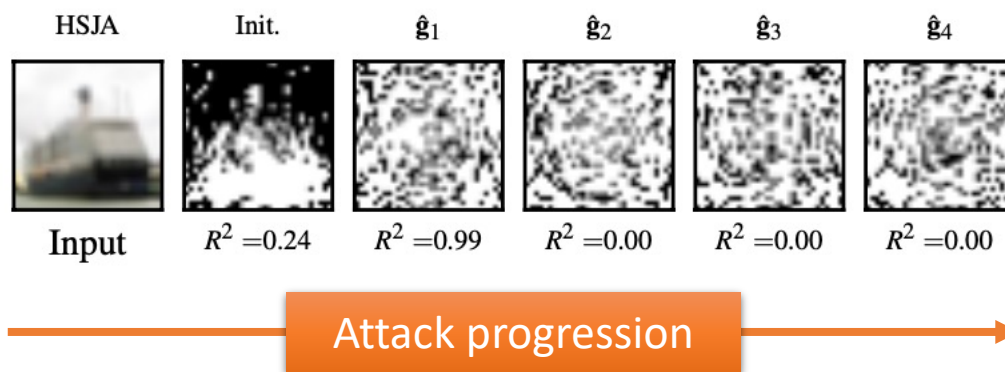     $\{R^2_{init}, R^2_{\hat{\mathbf{g}}_1}, \ldots, R^2_{\hat{\mathbf{g}}_n}\}$

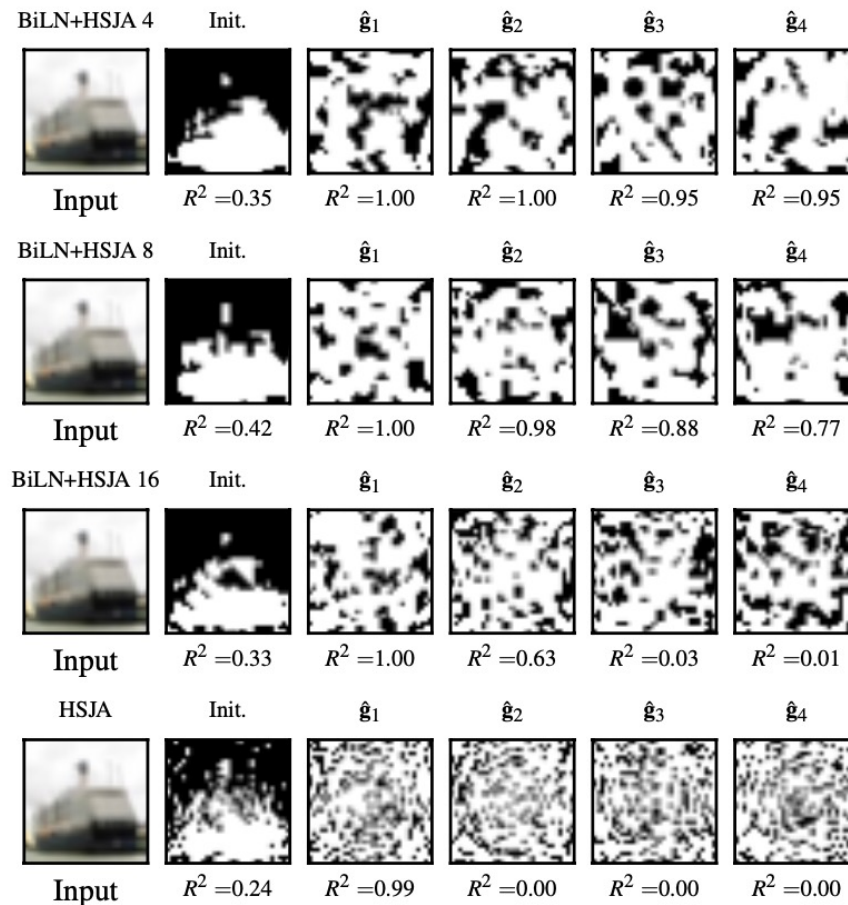Using (average) $R^2$ score of local models, we can answer the following:

1. Are hard-label queries sufficient to model the model's semantic structure in the query neighborhood?

2. Does dimension-reduction influence our structural knowledge?

A1. Yes:



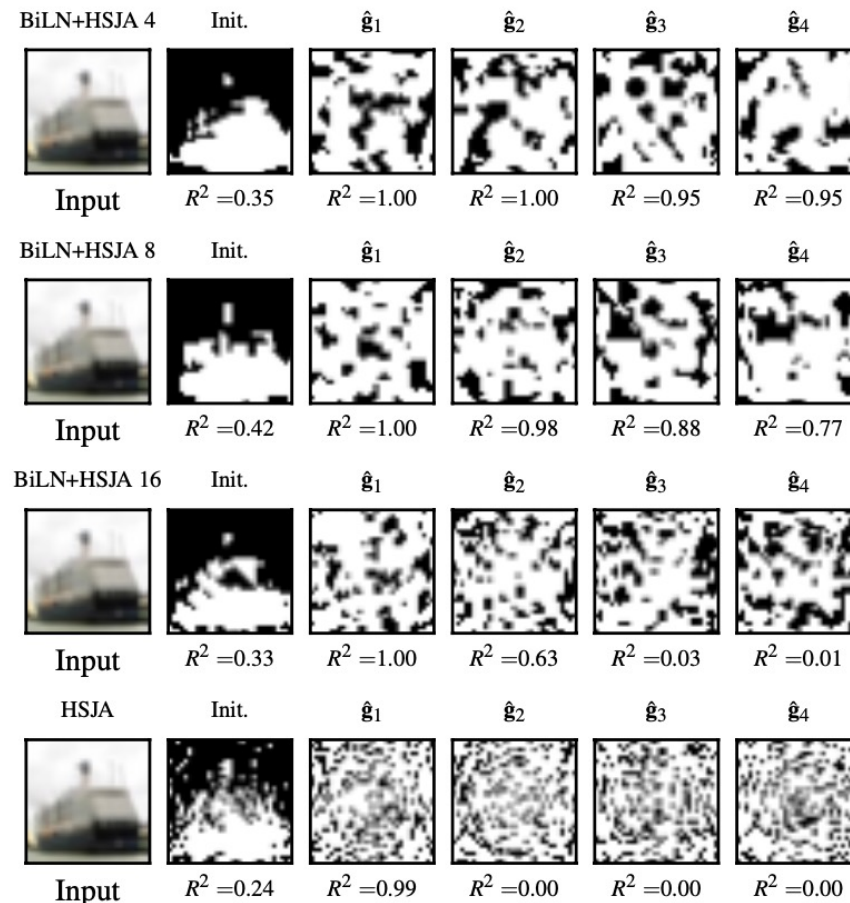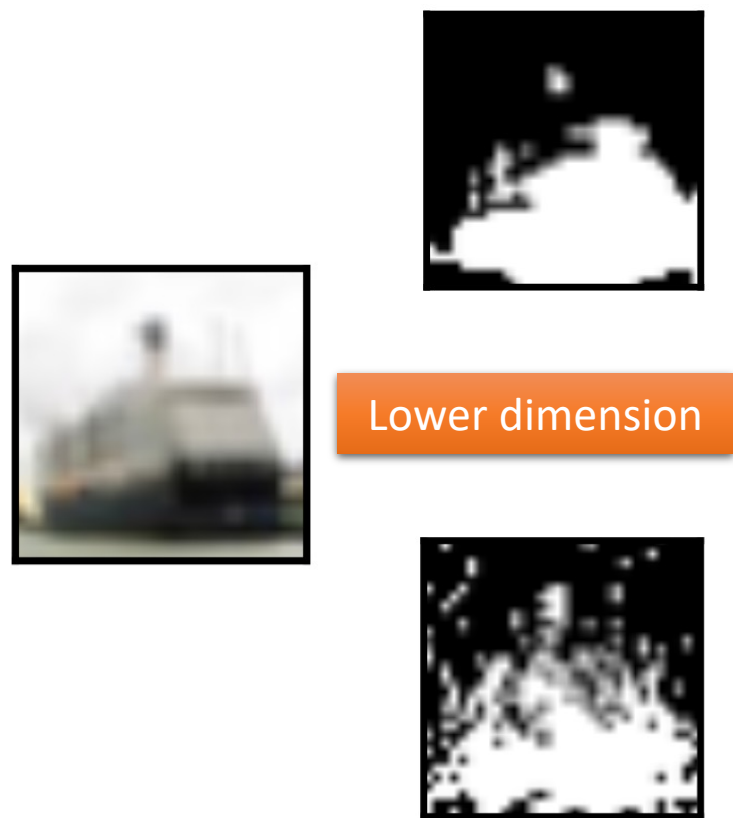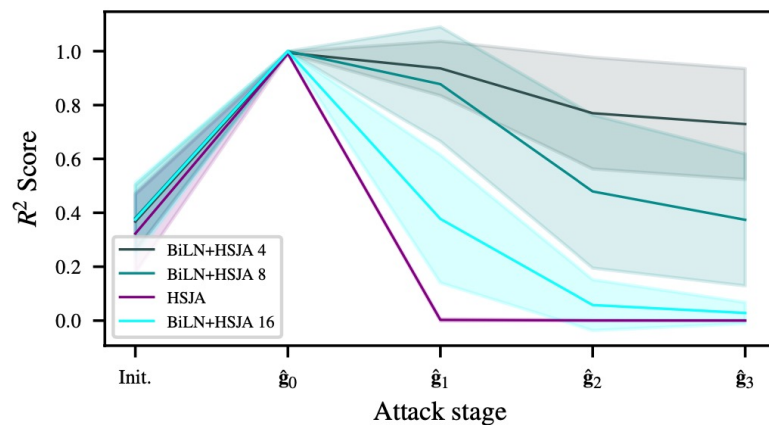| HSJA | Init. | $\hat{g}_1$ | $\hat{g}_2$ | $\hat{g}_3$ | $\hat{g}_4$ |
| Input | $R^2 = 0.24$ | $R^2 = 0.99$ | $R^2 = 0.00$ | $R^2 = 0.00$ | $R^2 = 0.00$ |

Attack progression

## A2. Dimension-reduction leads to finer-grained structural information:

Lower dimension

## A2. Dimension-reduction leads to finer-grained structural information:



Lower dimension

## A2. Dimension-reduction leads to finer-grained structural information:



| | Attack Variant | $\bar{R}^2$ | FID | SR@40k $(\epsilon=0.031)$ | LPIPS |
|---|---|---|---|---|---|
| Madry CIFAR-10 | HSJA | 0.259 | 0.244 | 0.272 | $0.676\pm0.275$ |
| | $\rightarrow$ BiLN 16 | 0.363 | 0.074 | **0.298** | $0.654\pm0.277$ |
| | $\rightarrow$ BiLN 8 | 0.624 | 0.026 | 0.224 | $0.668\pm0.304$ |
| | $\rightarrow$ BiLN 4 | **0.779** | **0.026** | 0.130 | $0.709\pm0.345$ |
| Natural CIFAR-10 | HSJA | 0.263 | 0.240 | **1.000** | $0.496\pm0.211$ |
| | $\rightarrow$ BiLN 16 | 0.368 | 0.085 | 0.984 | $0.543\pm0.227$ |
| | $\rightarrow$ BiLN 8 | 0.622 | 0.028 | 0.826 | $0.624\pm0.253$ |
| | $\rightarrow$ BiLN 4 | **0.759** | **0.012** | 0.472 | $0.651\pm0.297$ |

Second part of original question, what is the worst-case attack analysis?

Formulate adaptive attacks based on Algorithm 1 & 2, denoted MC and DynBiLN (cyan):

CIFAR-10

| Attack Variant | FID | SR AUC ($\epsilon=0.031$) |
|---|---|---|
| HSJA | 0.253 | 0.537 |
| $\rightarrow$ BiLN 4 | 0.026 ↓ | 0.342 |
| $\rightarrow$ BiLN 8 | 0.023 ↓ | 0.574 ↑ |
| $\rightarrow$ BiLN 16 | 0.074 ↓ | 0.720 ↑ |
| MC HSJA | 0.213 ↓ | 0.545 ↑ |
| $\rightarrow$ BiLN 4 | **0.022** ↓ | 0.356 |
| $\rightarrow$ BiLN 8 | 0.026 ↓ | 0.577 ↑ |
| $\rightarrow$ BiLN 16 | 0.068 ↓ | 0.705 ↑ |
| $\rightarrow$ DynBiLN | 0.030 ↓ | 0.607 ↑ |
| RayS | 0.057 | **1.000** |

ImageNet

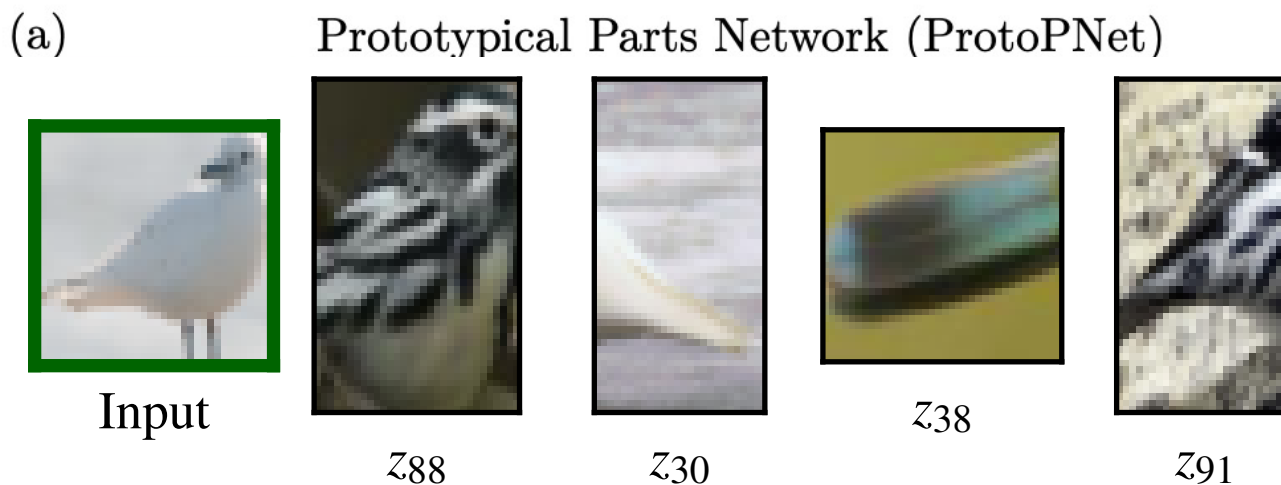| Attack Variant | FID | SR AUC ($\epsilon=0.031$) |
|---|---|---|
| HSJA | 1.541 | 0.344 |
| $\rightarrow$ BiLN 16 | 0.312 ↓ | 0.777 ↑ |
| $\rightarrow$ BiLN 32 | 1.085 ↓ | 0.771 ↑ |
| $\rightarrow$ BiLN 64 | 2.567 | 0.655 ↑ |
| MC HSJA | 1.591 | 0.331 |
| $\rightarrow$ BiLN 16 | **0.271** ↓ | 0.772 ↑ |
| $\rightarrow$ BiLN 32 | 1.079 ↓ | 0.771 ↑ |
| $\rightarrow$ BiLN 64 | 2.287 | 0.615 ↑ |
| $\rightarrow$ DynBiLN | 0.657 ↓ | 0.774 ↑ |
| RayS | 0.302 | **1.000** |

Two main questions:

1.  Is agent concept resolution possible from an *adversary's* perspective? What's the worst-case security analysis?

2.  It isn't always realistic to have a human-in-the-loop, or fine-grained data labels. What are ways around this?

A1: Expansion of our previous hard-label paper, in submission to IEEE SaTML 2023.

A2: We presented an initial idea @ NAACL 2022. Introduced expansion during previous meeting. Then worked on it as part of AFRL 2022 internship.

# First, agents learn human-interpretable perceptual knowledge priors:



(a) Prototypical Parts Network (ProtoPNet)

Input    $z_{88}$    $z_{30}$    $z_{38}$    $z_{91}$

C. Chen, O. Li, C. Tao, A. J. Barnett, J. Su, and C. Rudin, "This Looks Like That: Deep Learning for Interpretable Image Recognition," *arXiv:1806.10574 [cs, stat]*, 2019.

- Sender solves two joint tasks:
    1. Learn to embed their top-1 activate structure ($\mathbf{z}^S$) in the message
    2. Learn to describe the target objects

- Receiver solves two joint tasks:
    1. Learn to reconstruct the sender's top-1 structure ($rec(\mathbf{z}^S)$) from the message (*reconstruction loss*)

$$\mathcal{L}_{rec}(\mathbf{z}^S, rec(\mathbf{z}^S)) = \frac{1}{L} \sum_{l=1}^{L} |\mathbf{z}_{(l)}^S - rec(\mathbf{z}_{(l)}^S)|$$

    2. Learn to signal the correct target object (*classification loss*)

$$\mathcal{L}_{cls}(\mathbf{t}) = -\sum_{l=1}^{L} \alpha \log p(y_{(l)} = \mathbf{t} \mid msg_{(l)});$$

$$\boxed{\mathcal{L} = \mathcal{L}_{cls} + \mathcal{L}_{rec}}$$

Consider a gradual expansion of the sender agent's concept allowance, as in "Tatanka" clip from Dances with Wolves:
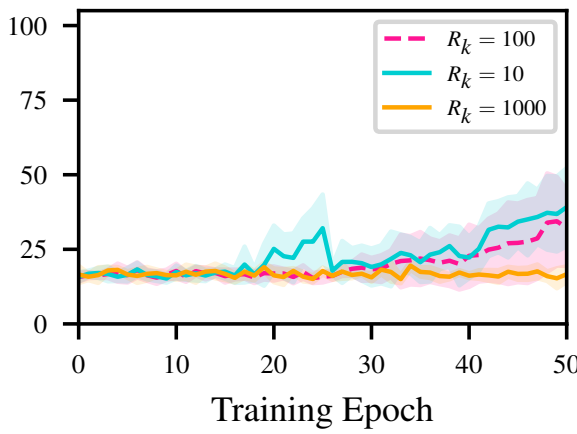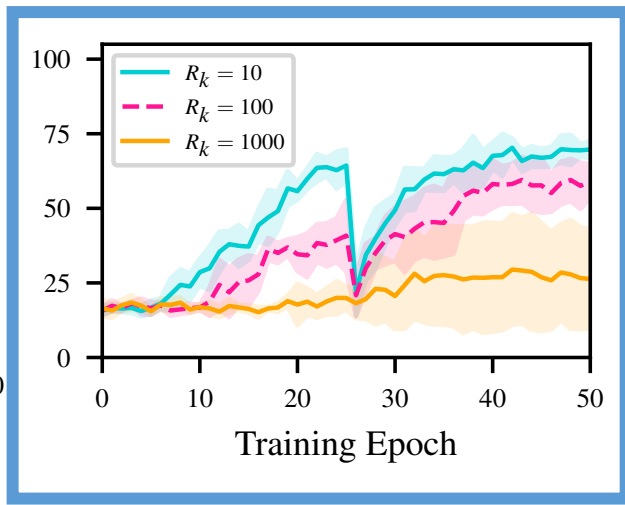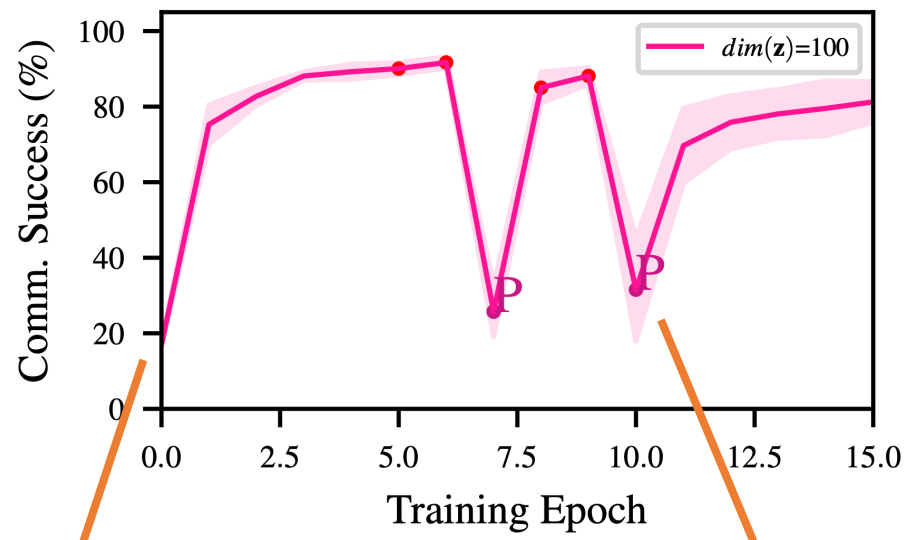
Baseline

Reification
(Ours)

Reification
"Tatanka" Variant
(Ours)

- Senders solve ~~two~~ three joint tasks:
  1. Learn to embed their top-1 activate structure in the message
  2. Learn to describe the target objects
  3. Update knowledge structure based on embedding difficulty

- Receivers solve ~~two~~ three joint tasks:
  1. Learn to reconstruct the sender's top-1 structure from the message
  2. Learn to signal the correct target object
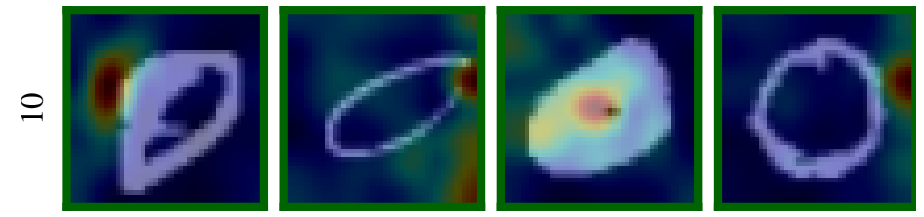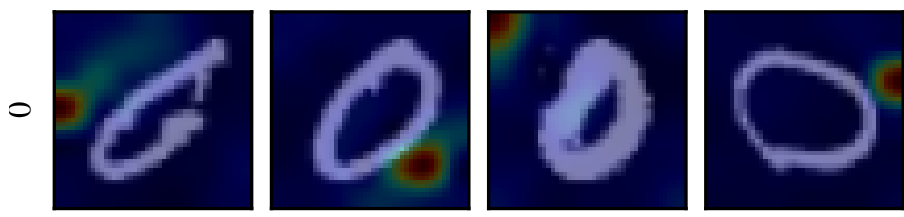  3. Update knowledge structure based on perceived utility of sender structure

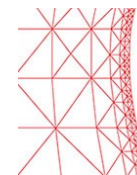# Training instability and automatic recovery:



P = human-aligned correction

Adversaries can learn semantic structure in a neighborhood around a sample, and this informs geometric interpretation of generalization errors.

- Can we get the global semantic structure with few samples and queries? Implication: leakage of learned manifold
- Connection to diffusion models

Semiotic learning offers an avenue for automatic structural validation, without explicit labels!

Graduation: Dec. 2022

Joining UDRI in January

w.garcia@ufl.edu