

# Probing Large Multimodal Models (LMMs) via Semantic Information Pursuit

IARPA BENGAL Proposer's Day

**Johns Hopkins University:** Rama Chellappa

**University of Pennsylvania:** René Vidal

# Recent DoD Projects

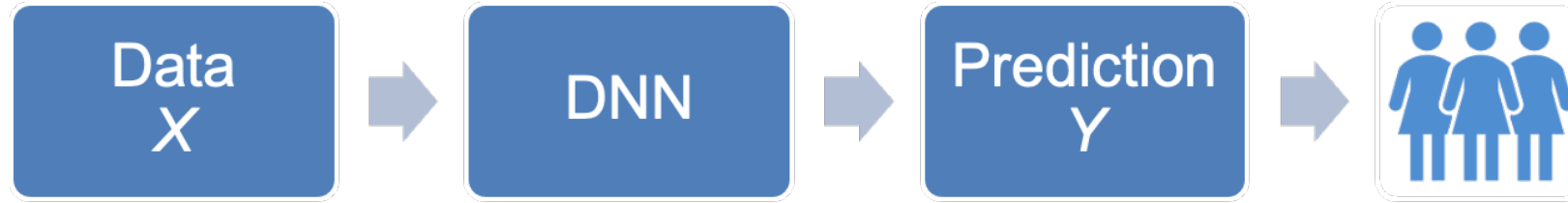
- **Rama Chellappa:**

- IARPA: WRIVA, BRIAR, DIVA, JANUS
- ONR MURI: Foundations of Deep Learning
- ARO MURI: Semantic Information
- DARPA: Guaranteeing AI Robustness Against Deception (GARD)

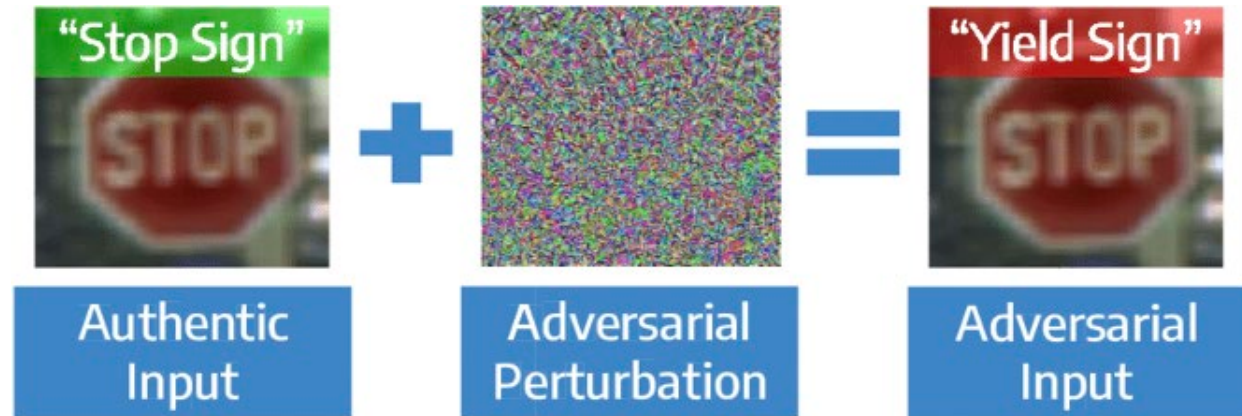
- **René Vidal:**

- IARPA: WRIVA, BRIAR, DIVA
- ONR MURI: Control and Learning Enabled Verifiable Robust AI (CLEVR-AI)
- ARO MURI: Semantic Information
- DARPA: Guaranteeing AI Robustness Against Deception (GARD)
- DARPA: Reverse Engineering Deception (RED)

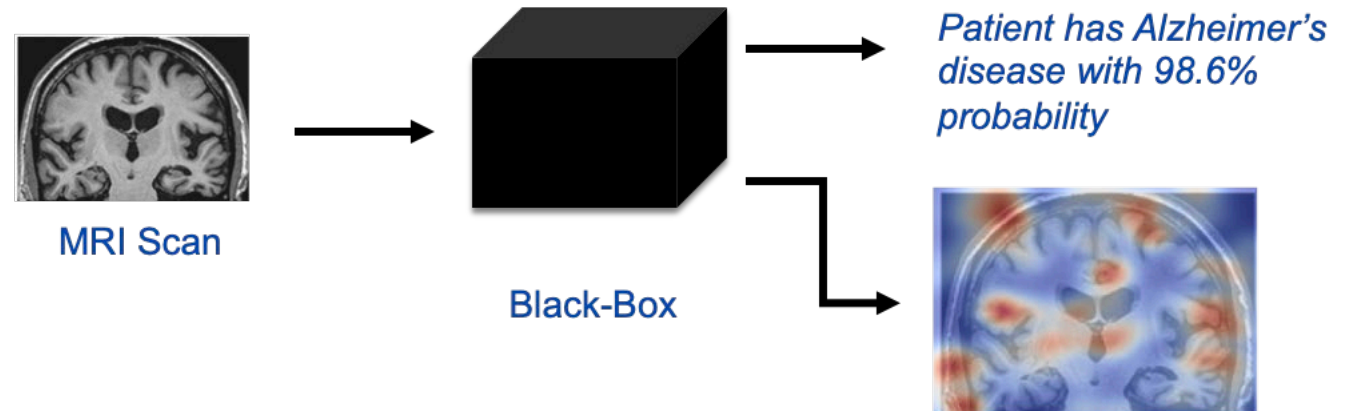
# Guarantees of Performance of AI Methods



- **Robustness:**



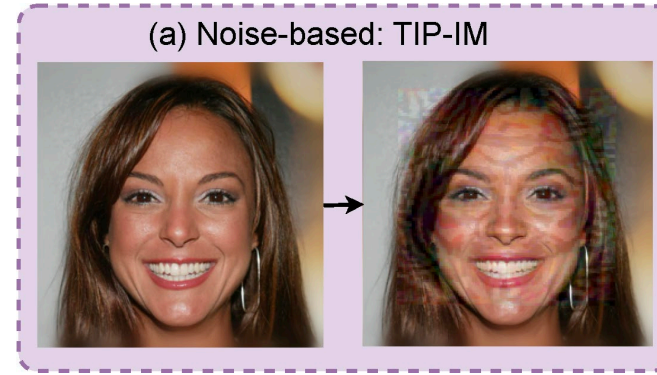
- **Explainability:**



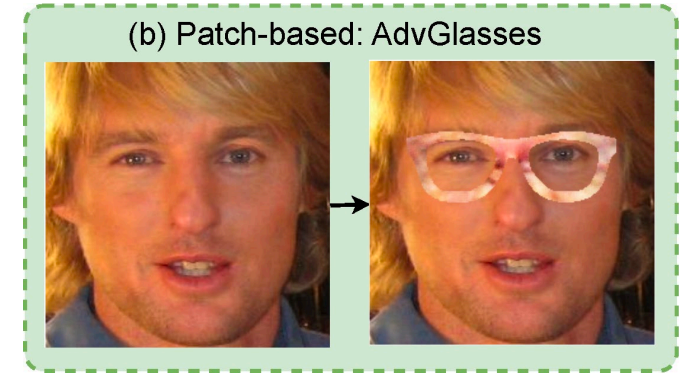
# Prior Work: Diffusion-Model Based Attacks

- **DiffProtect:**

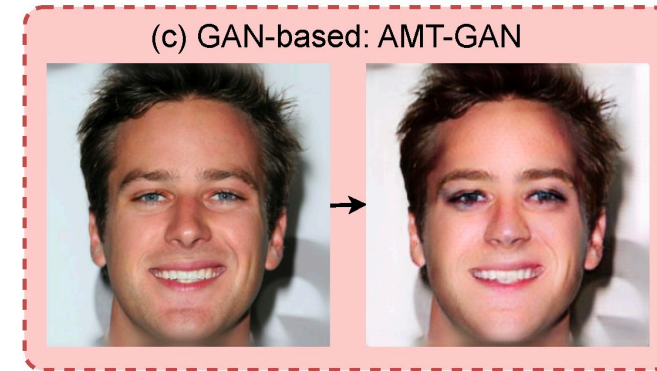
- Facial privacy protection
- Tradeoff between attack performance and visual quality



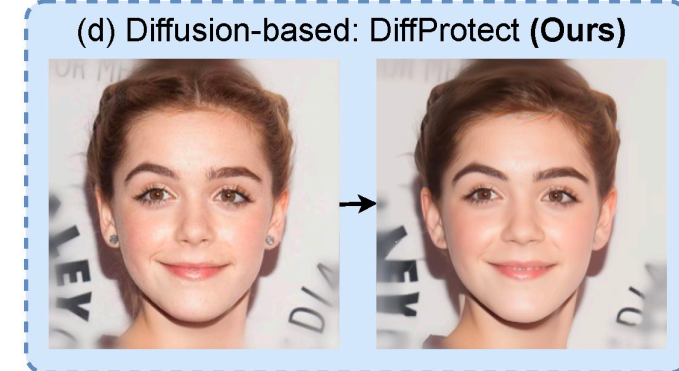
Visual Quality ↓  
Attack Performance ↑



Visual Quality ↓  
Attack Performance ↑



Visual Quality ↑  
Attack Performance ↓



Visual Quality ↑  
Attack Performance ↑

# Prior Work: Explainable AI by Design

- **Semantic Information Pursuit:**

- Map input to sequence of “questions” and “answers”
- Make predictions based on “most informative” questions and answers

Input image  $x^{\text{obs}}$



Ask a sequence of interpretable queries about  $x^{\text{obs}}$

$q_1$ .	Has shape perching-like?	<b>Yes</b>
$q_2$ .	Has bill shape all-purpose?	<b>Yes</b>
$q_3$ .	Has belly color yellow?	<b>Yes</b>
$q_4$ .	Has upperparts color yellow?	<b>No</b>
$q_5$ .	Has throat color yellow?	<b>No</b>
$q_6$ .	Has breast color black?	<b>Yes</b>
$q_7$ .	Has belly color olive?	<b>Yes</b>

Predicted bird species

Green Jay with  
99% probability



# Prior Work: Explainable AI by Design

- **Semantic Information Pursuit:**

- LLMs generate “imperfect questions” and LMMs generate “imperfect” answers
- Robustness is achieved by selecting the “most informative” questions and answers

Input image  $x^{\text{obs}}$



Ask a sequence of interpretable queries about  $x^{\text{obs}}$

$q_1$ .	Has shape perching-like?	Yes
$q_2$ .	Has bill shape all-purpose?	Yes
$q_3$ .	Has belly color yellow?	Yes
$q_4$ .	Has upperparts color yellow?	No
$q_5$ .	Has throat color yellow?	No
$q_6$ .	Has breast color black?	Yes
$q_7$ .	Has belly color olive?	Yes

Predicted bird species

Green Jay with  
99% probability

# Proposed Work for BENGAL

- **Probing LMMs via Semantic Information Pursuit**
- Semantic information pursuit provides explanations for both “correct” and “incorrect” predictions, thus allowing us to “diagnose” the model
- When the model makes a mistake, we know why, and we can use the explanation to “correct mistakes”
- We can thus use semantic information pursuit to detect, characterize and mitigate LMM threats and vulnerabilities