# ABYSS:

## AI MEETS THE DARK SIDE

SPEAKER: STEVEN TIRTADJAJA

MCS 11

# TABLE OF CONTENTS

SPEAKER: STEVEN TIRTADJAJA

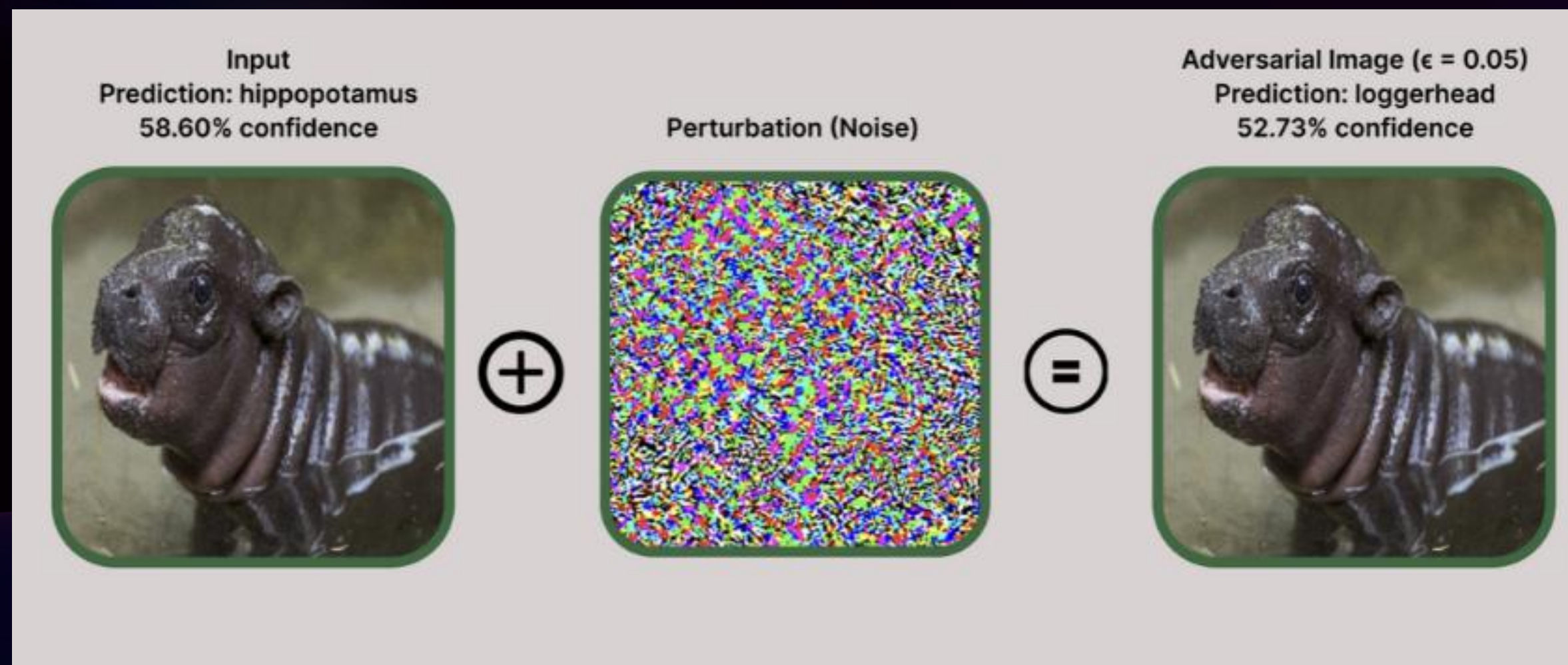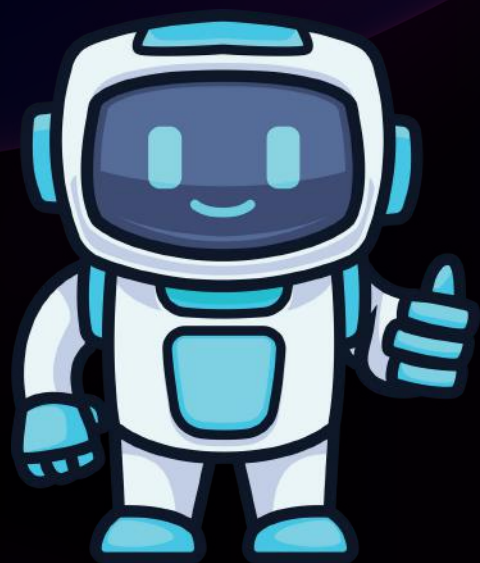# INTRODUCTION

- Rapid AI advancements

- Vulnerabilities in Image Classification AI Models

- Fast Gradient Sign Method
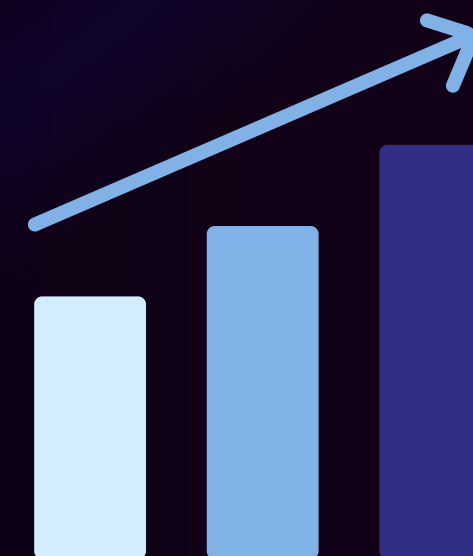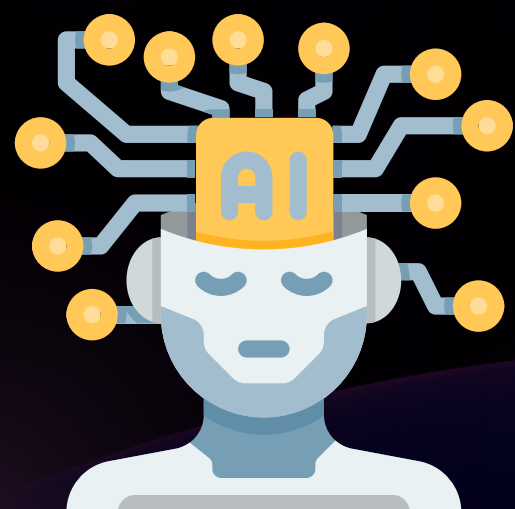
SPEAKER: STEVEN TIRTADJAJA

# FAST GRADIENT SIGN METHOD

# PROBLEM STATEMENT

Rapid advancements of AI have uncovered many vulnerabilities especially in image classification AI models

SPEAKER: STEVEN TIRTADJAJA
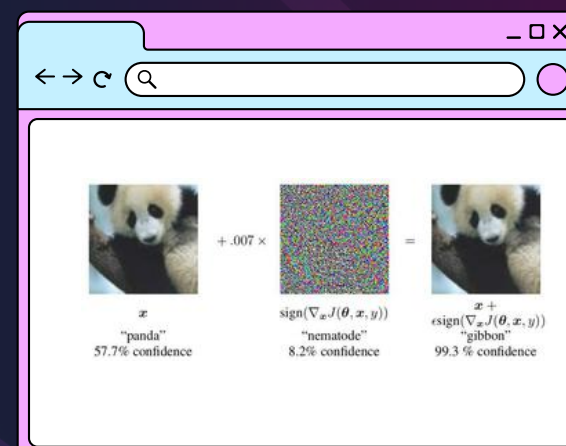
# PROJECT GOALS AND DELIVERABLES

- Examine the vulnerabilities of image classification AI models.
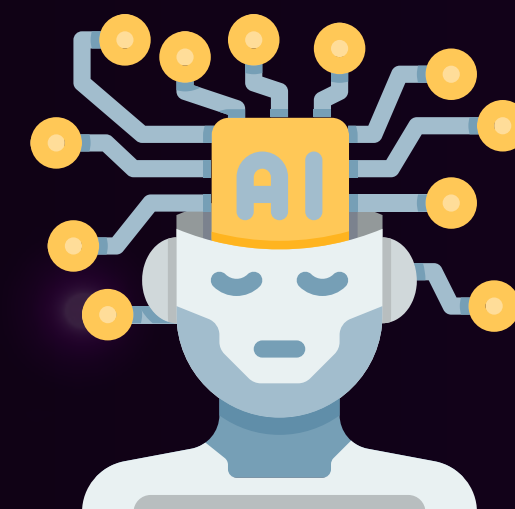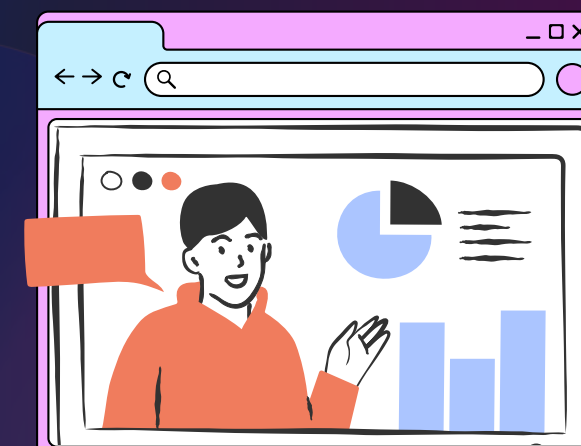- Raise awareness of AI weaknesses
- Demonstrate FGSM effectiveness

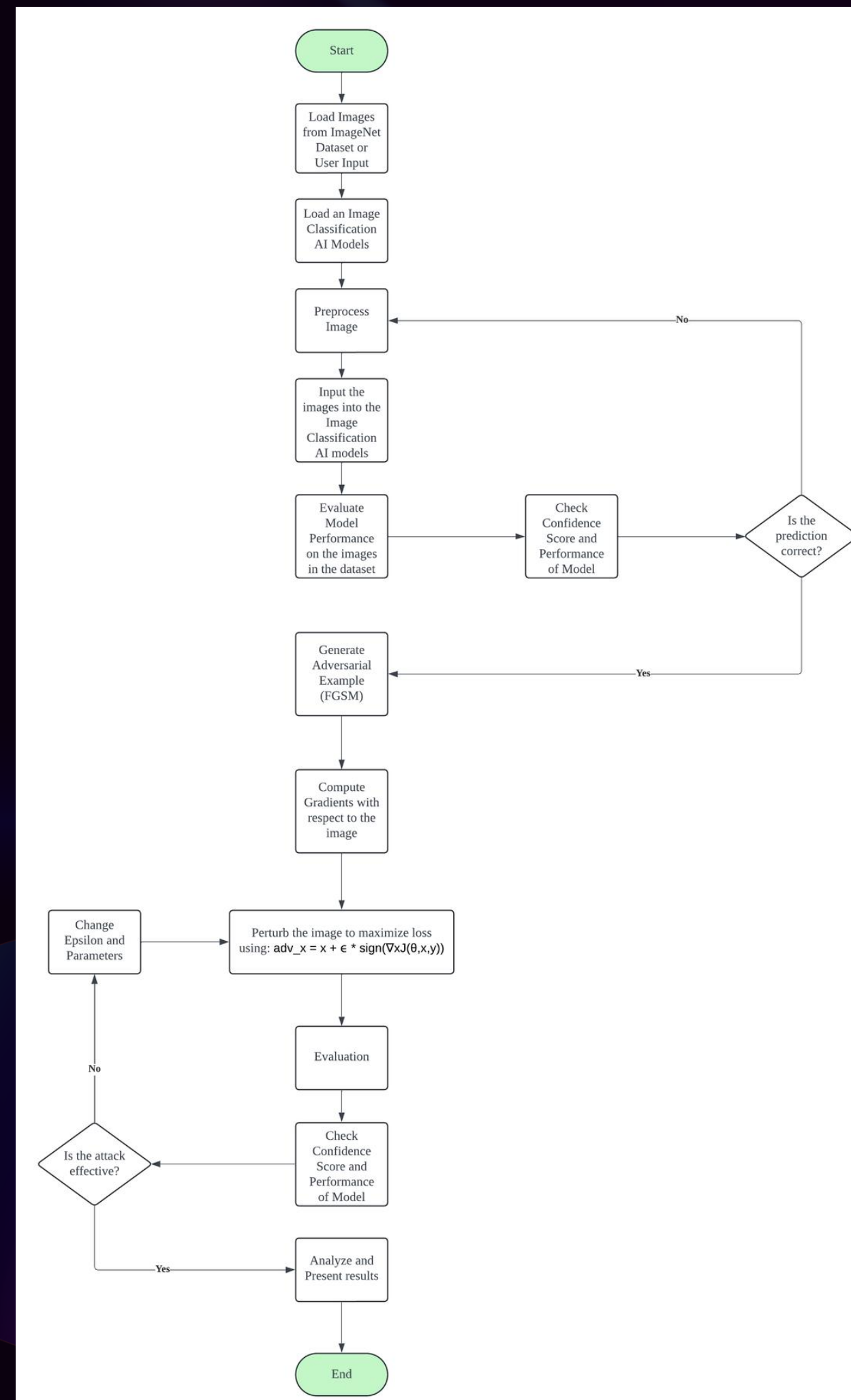**A functional web application to apply FGSM attacks**

**Visual side by side comparisons of original and adversarial images.**
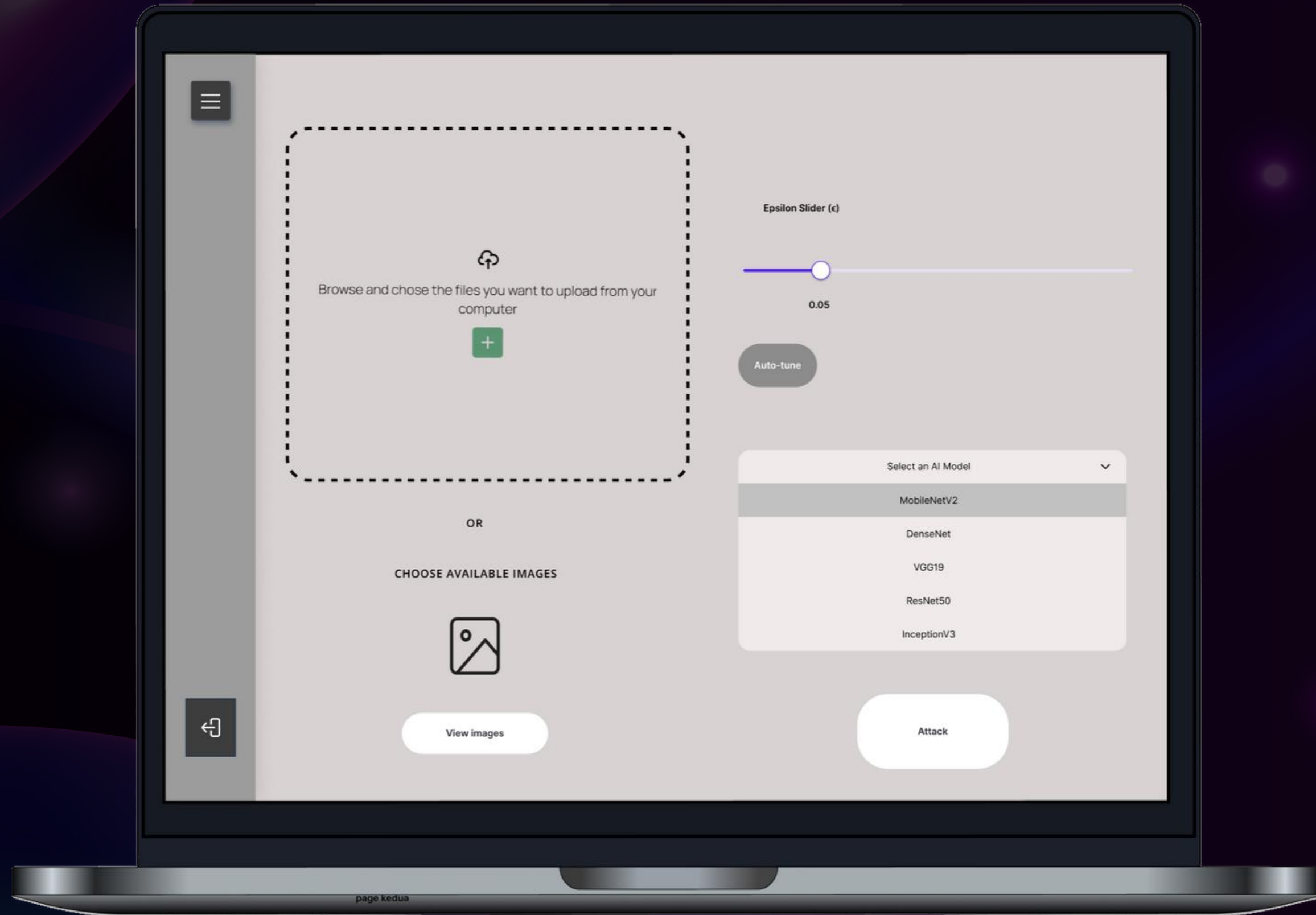
**Evaluate the impact of the adversarial attack**

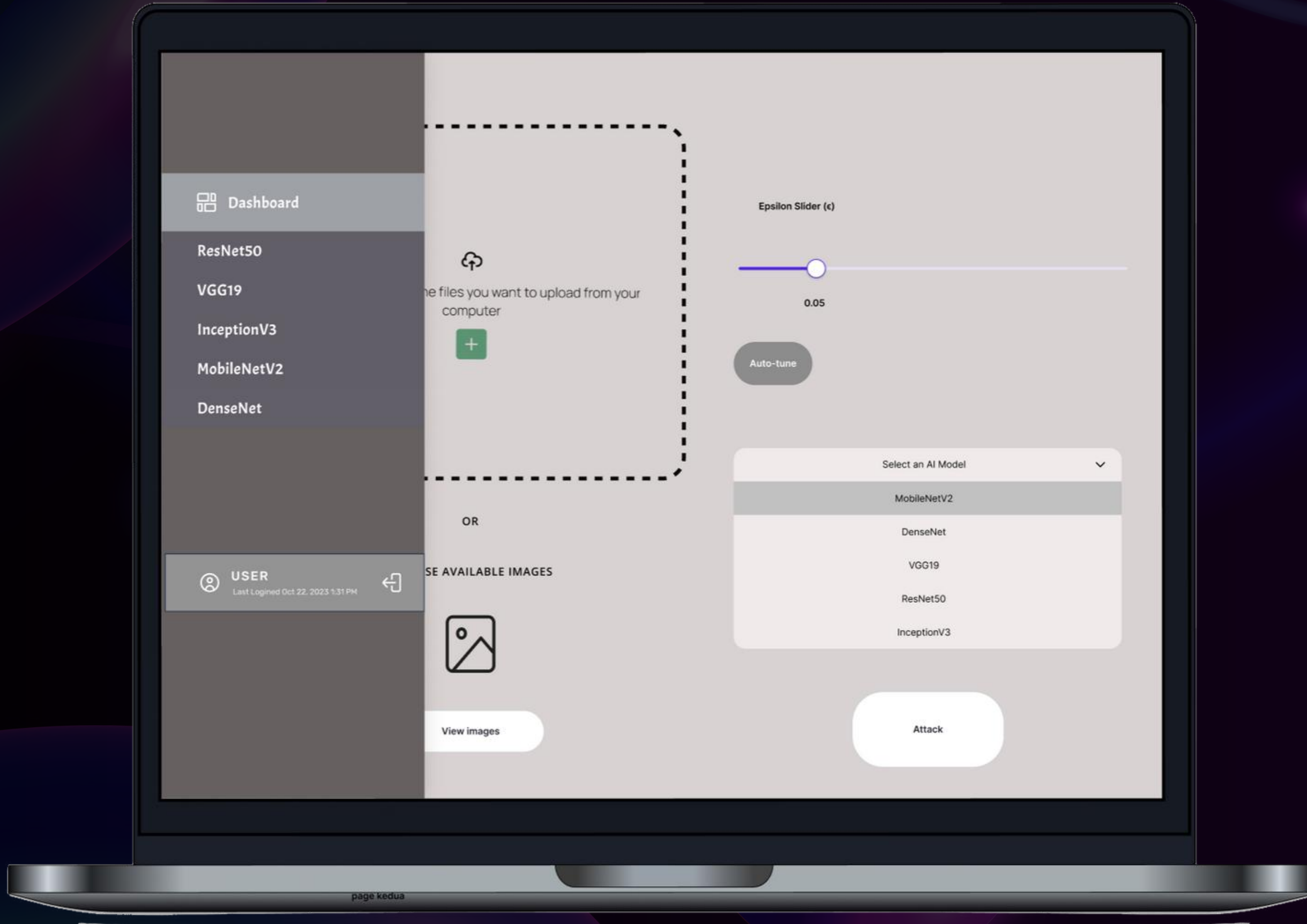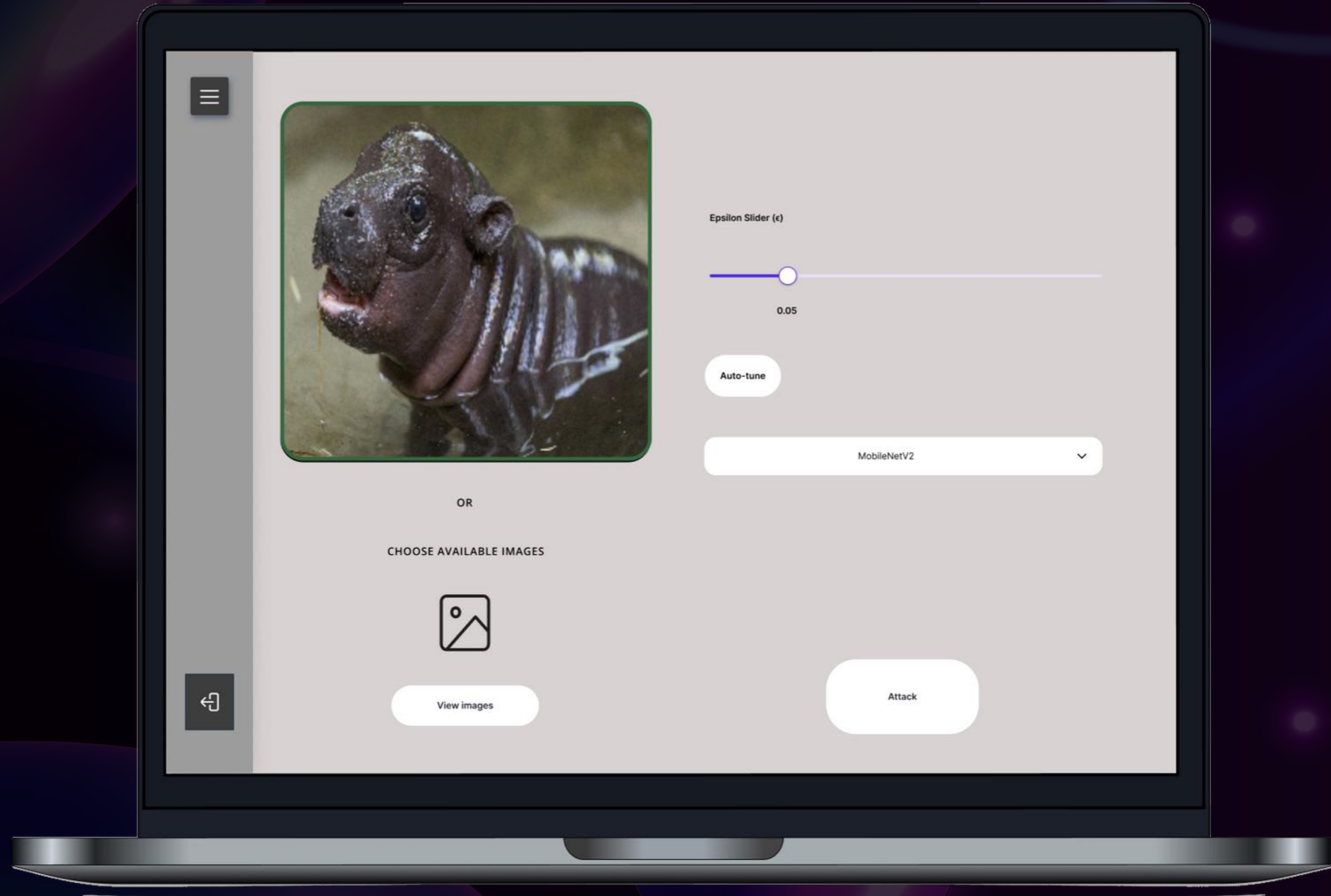SPEAKER: STEVEN TIRTADJAJA

# REPRESENTATIONS: FLOW CHART

SPEAKER: SHENG XIN CHUA

# REPRESENTATIONS: USER INTERFACE

SPEAKER: SHENG XIN CHUA

SPEAKER: SHENG XIN CHUA

SPEAKER: SHENG XIN CHUA

Result

Input
Prediction: hippopotamus
58.60% confidence

Perturbation (Noise)

Adversarial Image (ε = 0.05)
Prediction: loggerhead
52.73% confidence

SPEAKER: SHENG XIN CHUA

# PROJECT MANAGEMENT

## Google Drive

- Cloud storage and accessibility
- Large storage capacity
- Real time collaboration

## Github

- Track changes
- Repository management
- Ease of collaboration and integration with Visual Studio Code

SPEAKER: ALVIN ANDREAN

# KANBAN BOARD

## WHY?

- Improved communication and collaboration
- Clear workflow visualization
- Supports Incremental Progress

SPEAKER: MARCHOLAS TJANGNAKA

# PLANS

| Literature Review on FGSM | Applying Cleverhans' FGSM on more pre-trained models | Experiment on different image datasets | Start Implementing our own FGSM Attack |

SPEAKER: MARCHOLAS TJANGNAKA

# THANK YOU!

MCS11