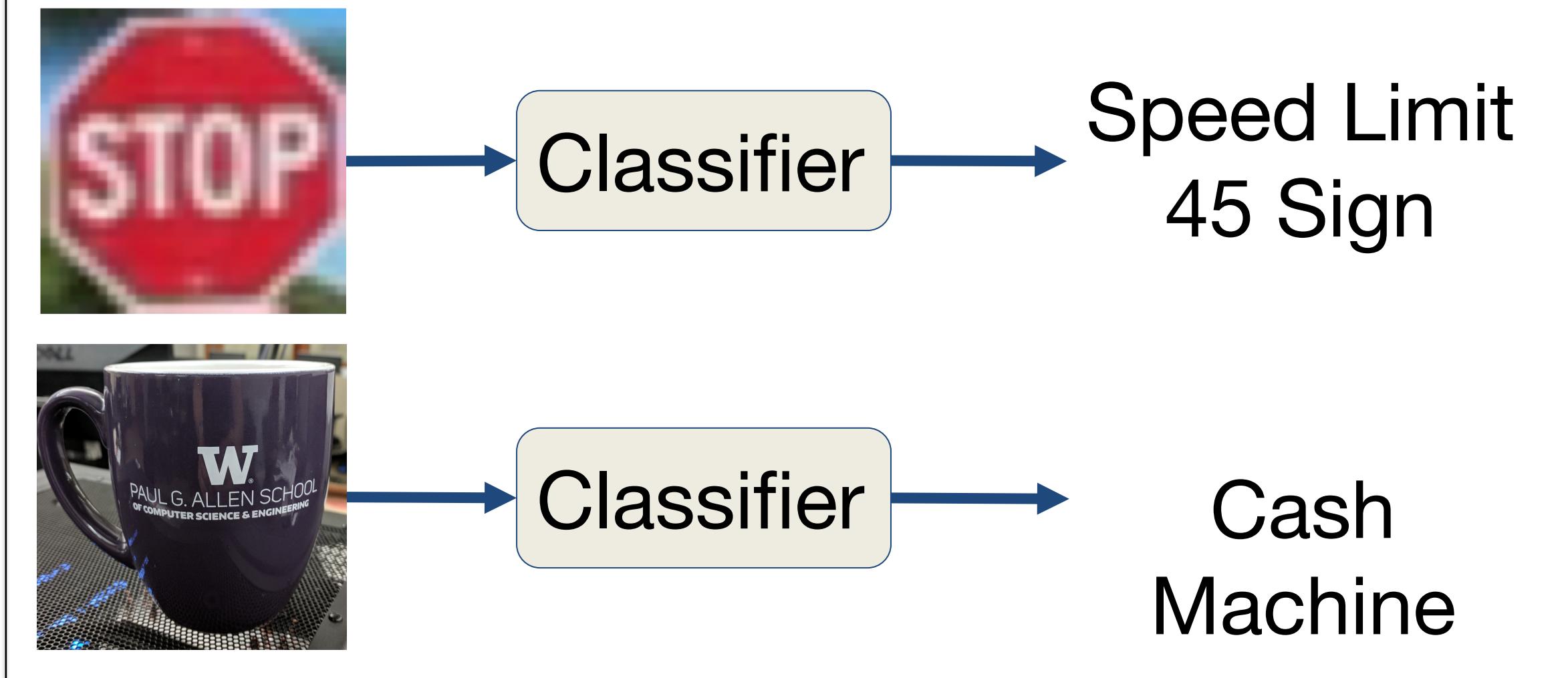


Robust Physical-World Attacks on Deep Learning Visual Classification

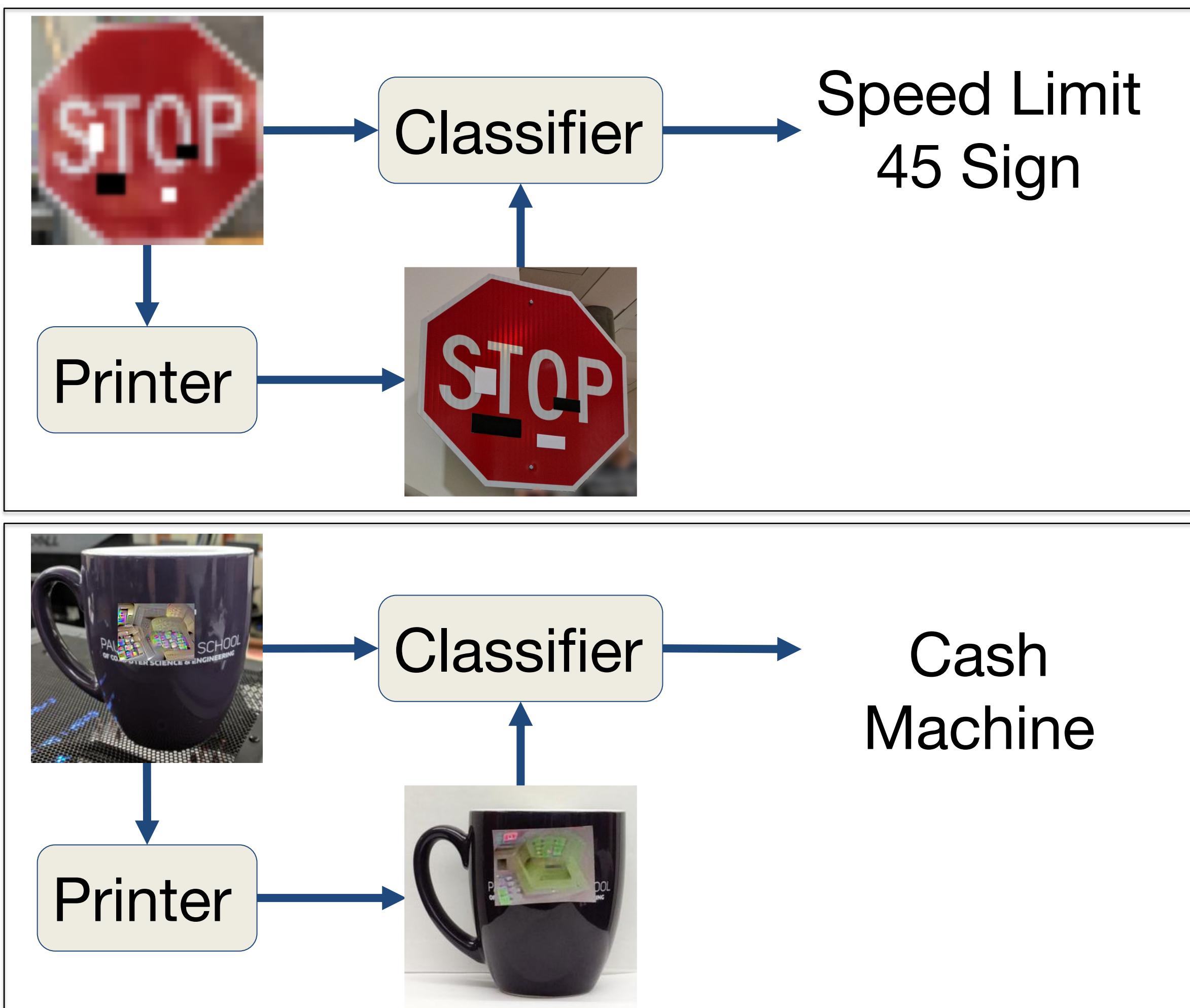
Kevin Eykholt¹, Ivan Etimov², Earlene Fernandes², Bo Li³, Amir Rahmati⁴, Chaowei Xiao¹, Atul Prakash¹, Tadayoshi Kohno², Dawn Song³
¹University of Michigan ²University of Washington ³University of California, Berkeley ⁴Samsung Research America and Stony Brook University

The Problem

Machine learning is vulnerable to **adversarial examples**



Can we create physical adversarial examples?



The Challenges of the Physical World

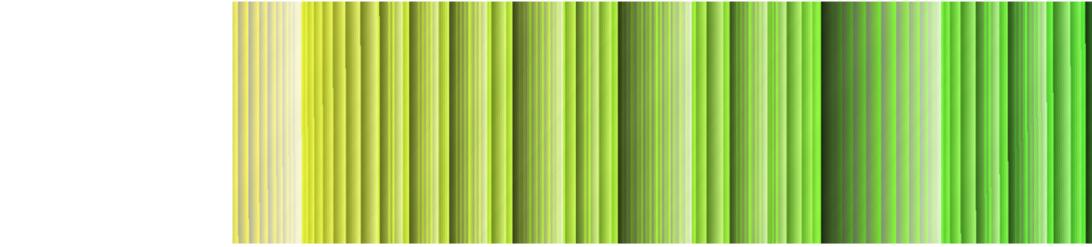
1) Viewer and Environmental Conditions



2) Spatial Constraints

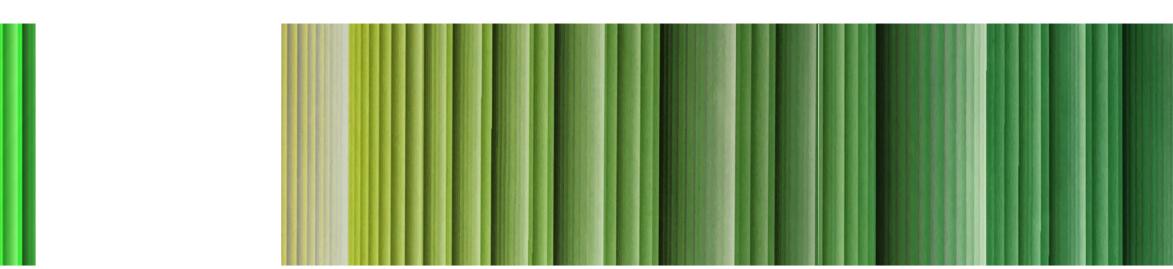
4) Fabrication Error

Digital Picture

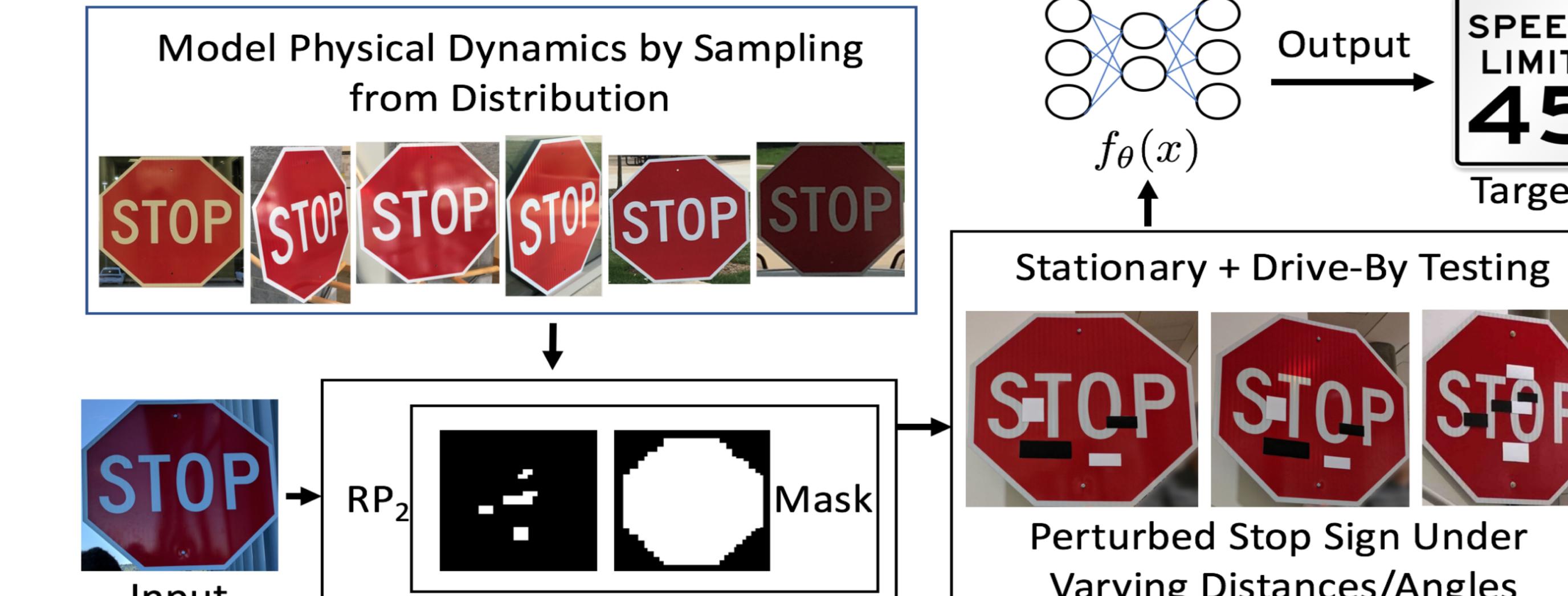


3) Limits of Imperceptibility

Printer Result



Robust Physical Perturbation – RP₂



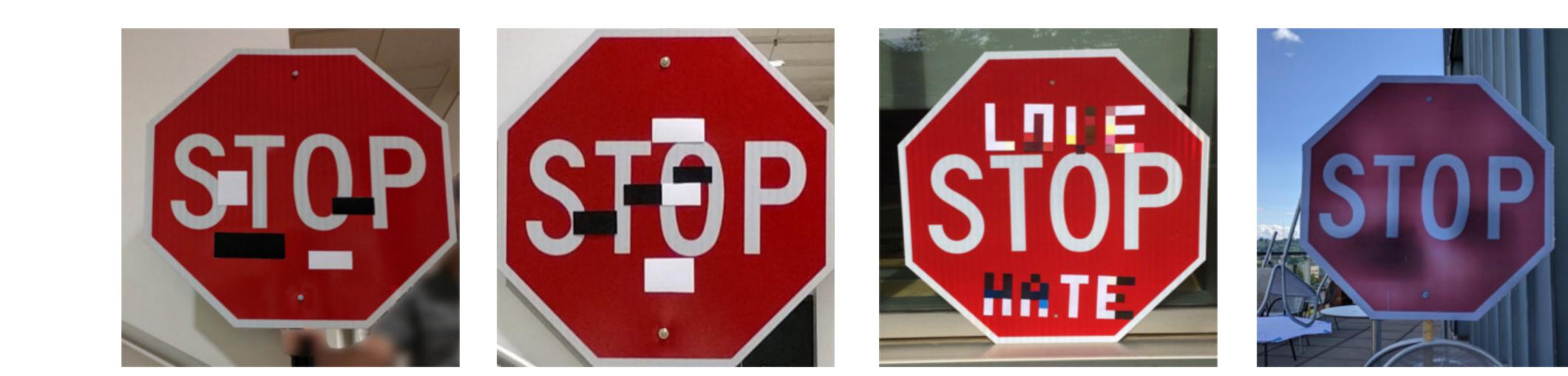
$$\text{Objective Function: } \underset{\delta}{\operatorname{argmin}} \lambda \| M_x \cdot \delta \| + NPS \\ + \mathbb{E}_{x_i \sim X^V} J(f_\theta(x_i + T_i(M_x \cdot \delta)), y^*)$$

Experiments

1) Vary Viewing Distances and Angles



2) Vary Mask Shapes



3) Vary Classification Model and Victim Object



| Attack Type | Success Rate | Attack Type | Success Rate |
|-------------|--------------|-------------|--------------|
| Camo Art | 100% | Right Turn | 73.33% |
| Camo Art v2 | 80% | Microwave | 90% |
| Graffiti | 66.67% | Cup | 71.4% |
| Subtle | 100% | | |

Visit <https://iotsecurity.eecs.umich.edu/#roadsigns> for samples images, videos, and other resources