



Oregon State
University

DARTS: Deceiving Autonomous Cars with Toxic Signs

Submitted to ACM CCS 2018 , *18 Feb 2018*

Chawin Sitawarin, Arjun Nitin Bhagoji

Shahab Nikkhoo



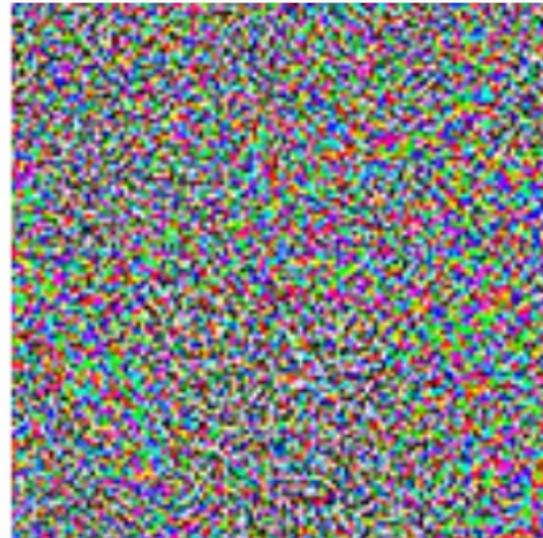
Evasion attack



“panda”

57.7% confidence

+ ϵ



=

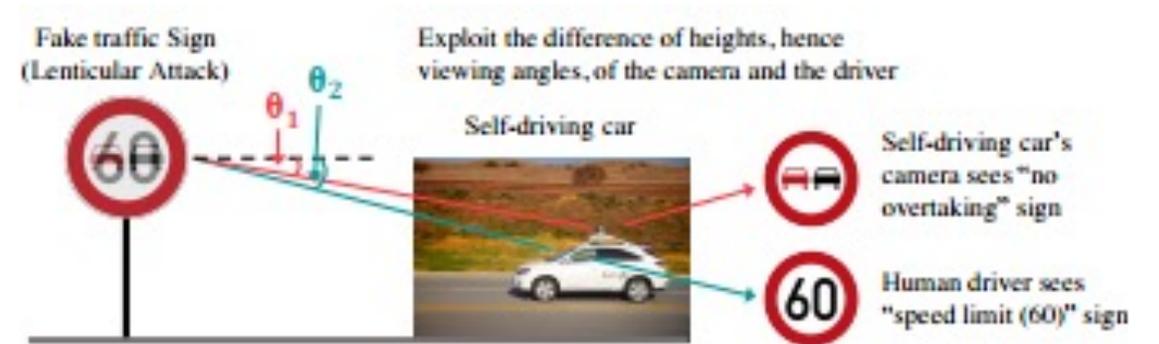
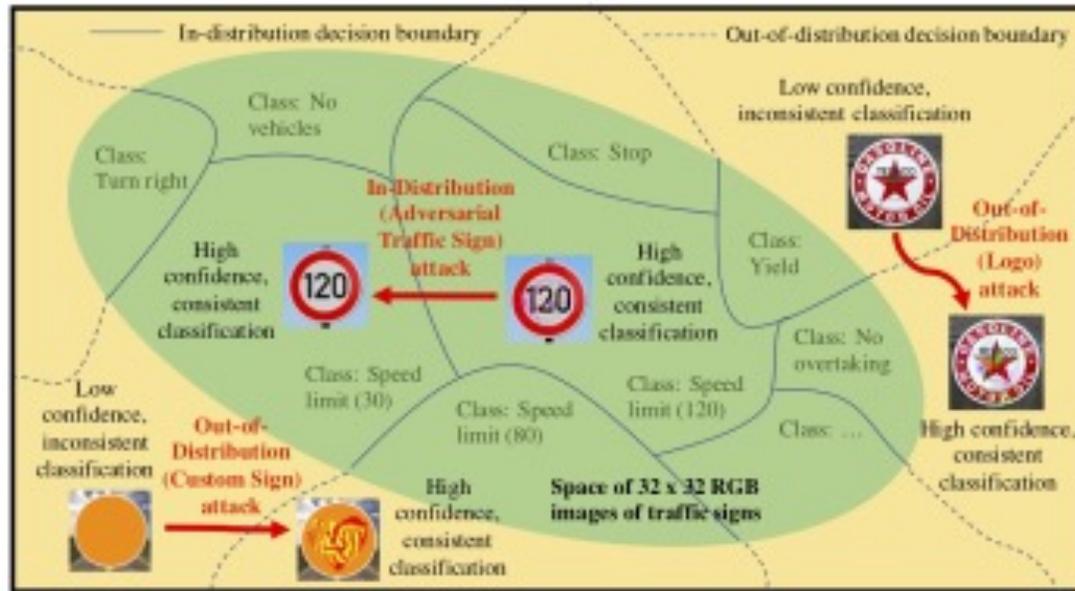


“gibbon”

99.3% confidence



Attack methods





attack methods

- In distribution
- Out of distribution
- Lenticular

Out-of-distribution Attacks

	Logo Attacks	Custom Sign Attacks
Original	 	 
Adversarial	 	 
Classified as:	Stop No overtaking	Speed limit (30) Stop

In-distribution Attacks

	Adversarial Traffic Signs
Original	 
Adversarial	 
Classified as:	Stop Speed limit (30)

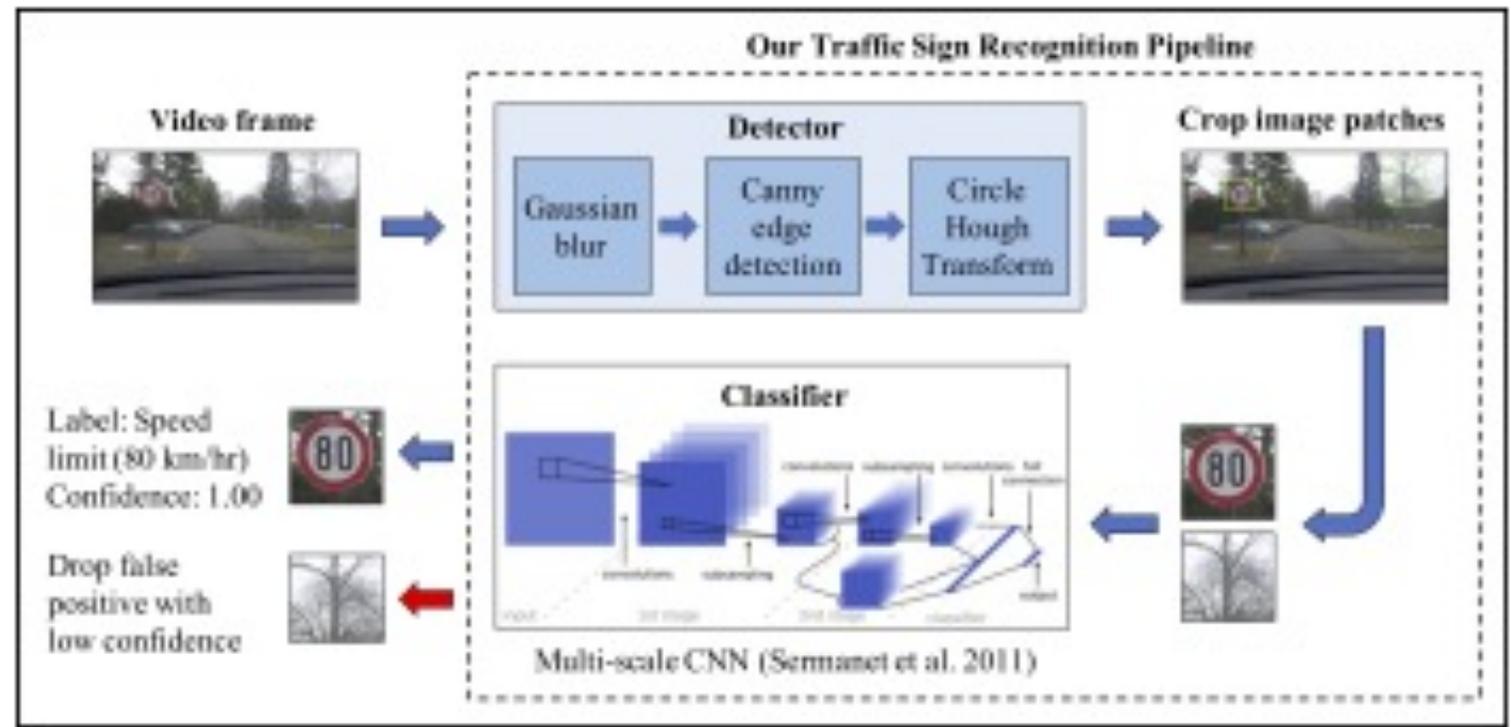
Lenticular Attacks

	Traffic sign – lenticular image	Logo – lenticular image
Straight view	 	
Angled view	 	
	Traffic sign – lenticular image	Logo – lenticular image



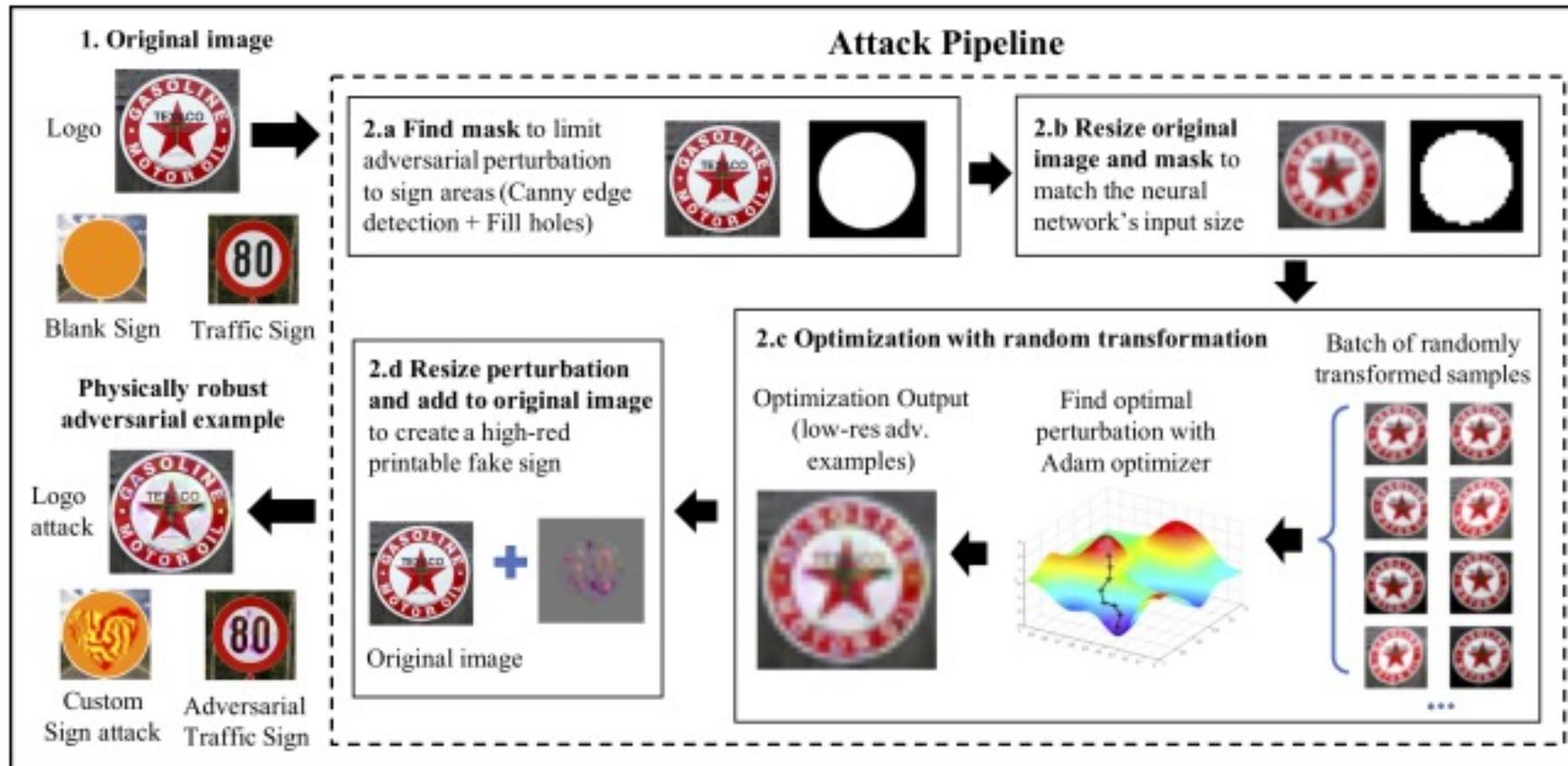
Classification pipeline

- Detection
- classification





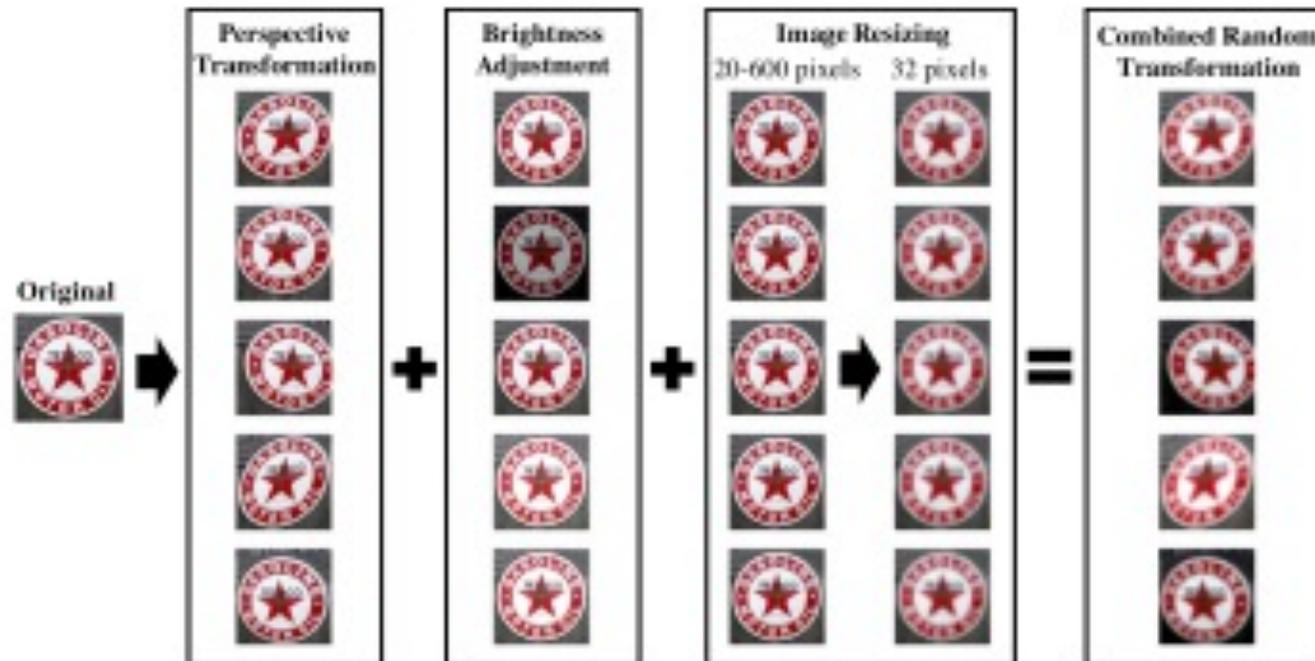
Attack pipeline





How to make it real?

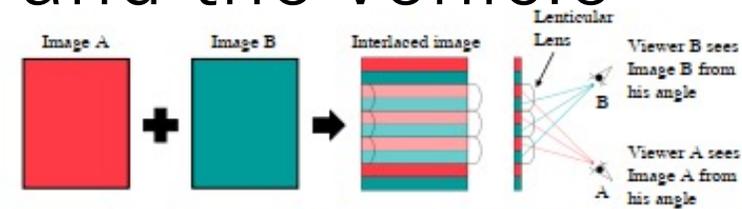
- Generate physical robust adversarial examples



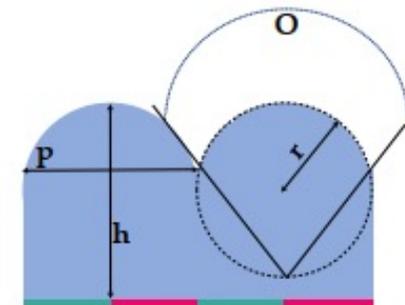


Lenticular attack

- Lenticular Printing attack is motivated by the key insight that the human driver and the vehicle mounted camera observe



(a) Illustration of the process of generating a lenticular image and its angle-dependent appearance.



(b) Several parameters (p , r , and h) determine the full angle of observation (O)



results

Attacks	Virtual attack success (VAS)	Simulated physical attack success (SPAS)	Avg. norm (L_1)	Avg. confidence
In-Distribution (auxiliary traffic data)	54.34 %	36.65 %	37.71	0.9721
Out-of-Distribution (Logo)	85.71%	65.07%	34.89	0.9753
Out-of-Distribution (Custom Sign)	29.44%	18.72%	N.A.	0.9508

Attacks	White-box (Avg. confidence)	Black-box (Avg. confidence)
In-Distribution (auxiliary traffic data)	92.82% (0.9632)	96.68% (0.9256)
Out-of-Distribution (Logo)	52.50% (0.9524)	32.73% (0.9172)
Out-of-Distribution (Custom Sign)	96.51% (0.9476)	97.71% (0.9161)



Oregon State University
College of Engineering

What about defense?

Adversarial training cannot defend against
Lenticular Printing and Out-of-Distribution



Discussion points

- They are using their own classifier for detection and recognizing the signs.
- Where are the logos ? Are they looks like the signs?
- This method is useful for AR?