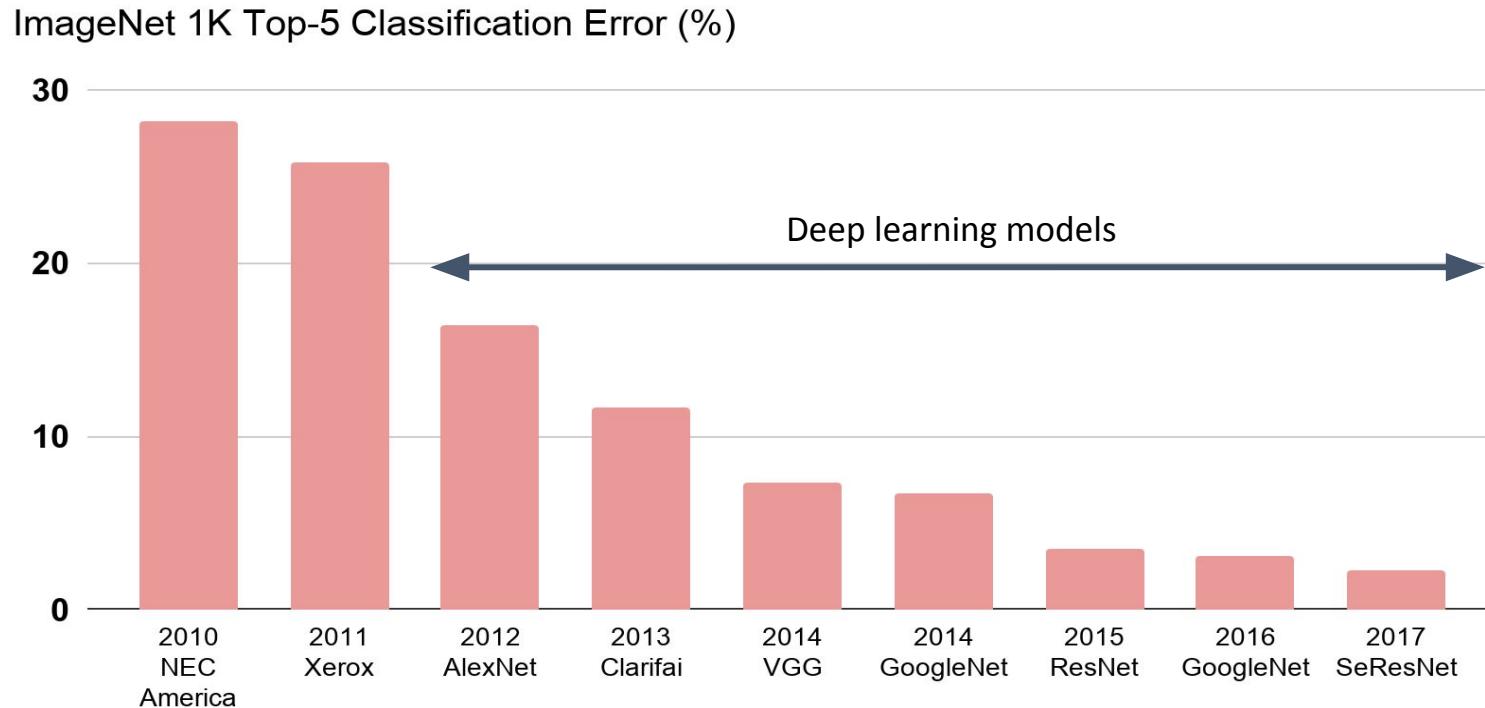


# Impact of Adversarial Examples on Deep Learning Models for Biomedical Segmentation

Utku Ozbulak, Arnout Van Messem, and Wesley De Neve

# Impact of Deep Learning Models on Computer Vision Problems

Deep learning methods drastically improved the state-of-the-art results in computer vision problems.



# Impact of Deep Learning Models on Machine Learning Problems

After their success on ImageNet and other datasets, the deep learning models are, at an increasing rate, being adopted on solving various machine learning problems such as:

- Facial recognition problems (e.g., personal identification)
- Problems related to Self-driving cars (e.g., lane detection)
- Problems related to smart-housing (e.g., voice commands)
- **Medical imaging problems (e.g., tumor detection)**

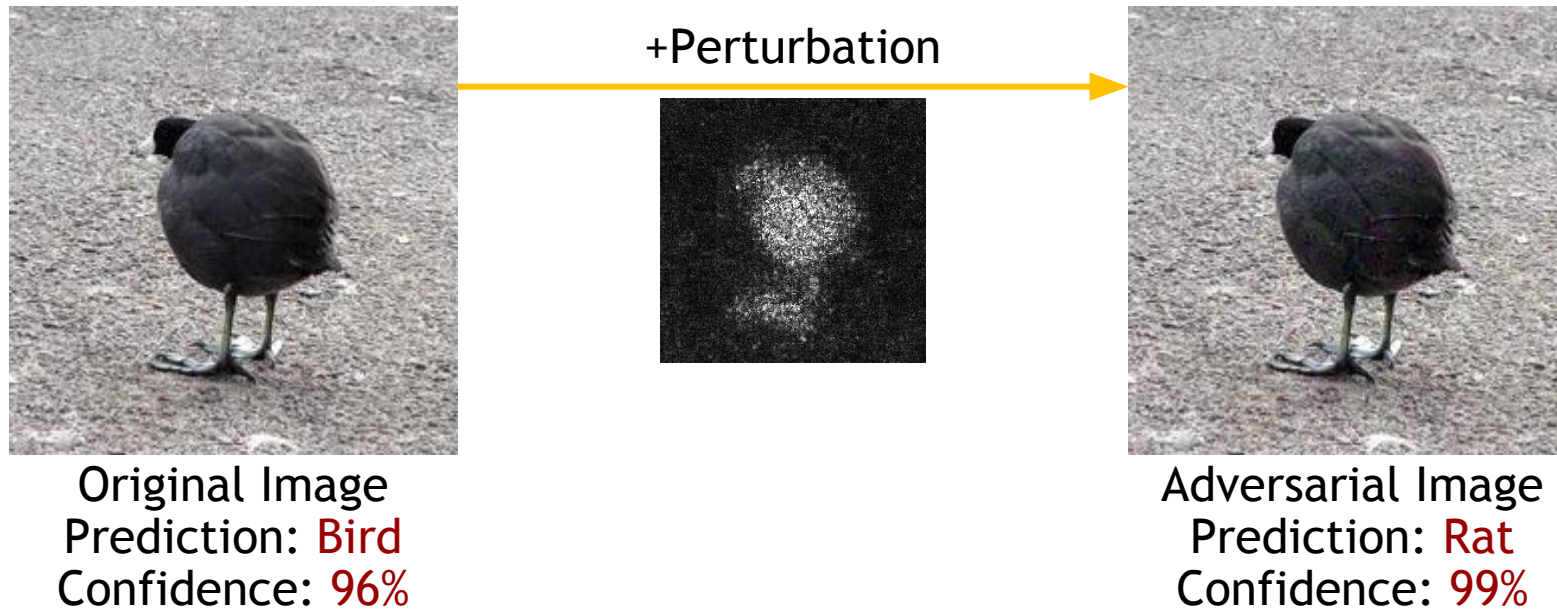
## Problems with Deep Learning Models

Although deep learning models are praised for their results on accuracy for complex problems, they are not perfect. A number of problems present in those models can be summarized as follows:

- Computational cost of training a model
- Reproducibility problems related to randomness
- Interpretability
- **Adversarial examples**

# What are Adversarial Examples?

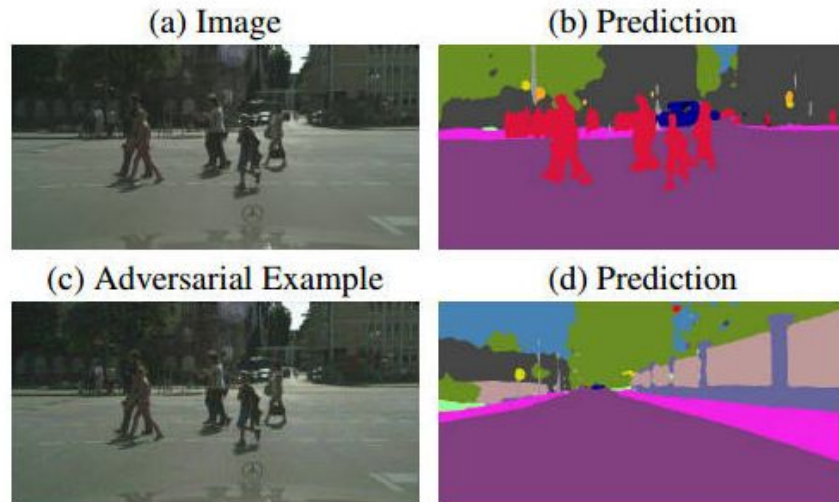
Adversarial examples are carefully crafted data points which force machine learning models into misclassification during testing phase. These malicious samples are often undetectable by humans.



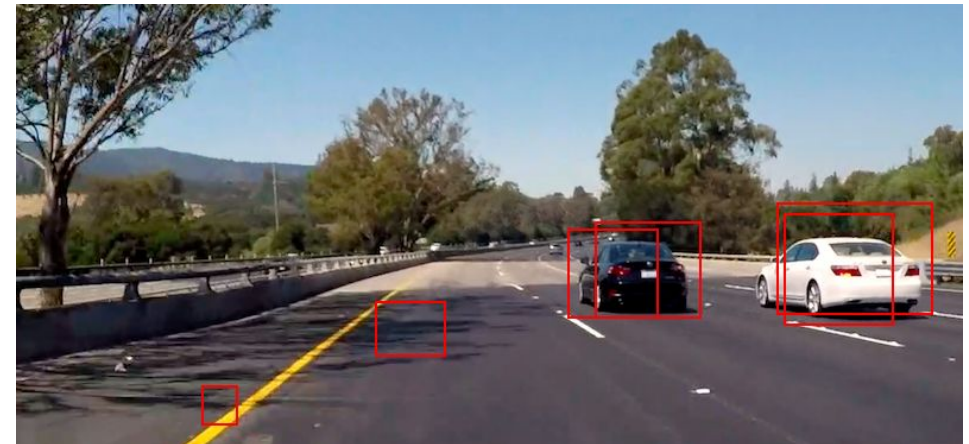
# Real World Consequences of Adversarial Examples

Adversarial examples with malicious intent reduce trust in automated systems (self-driving cars).

## Roadway Segmentation



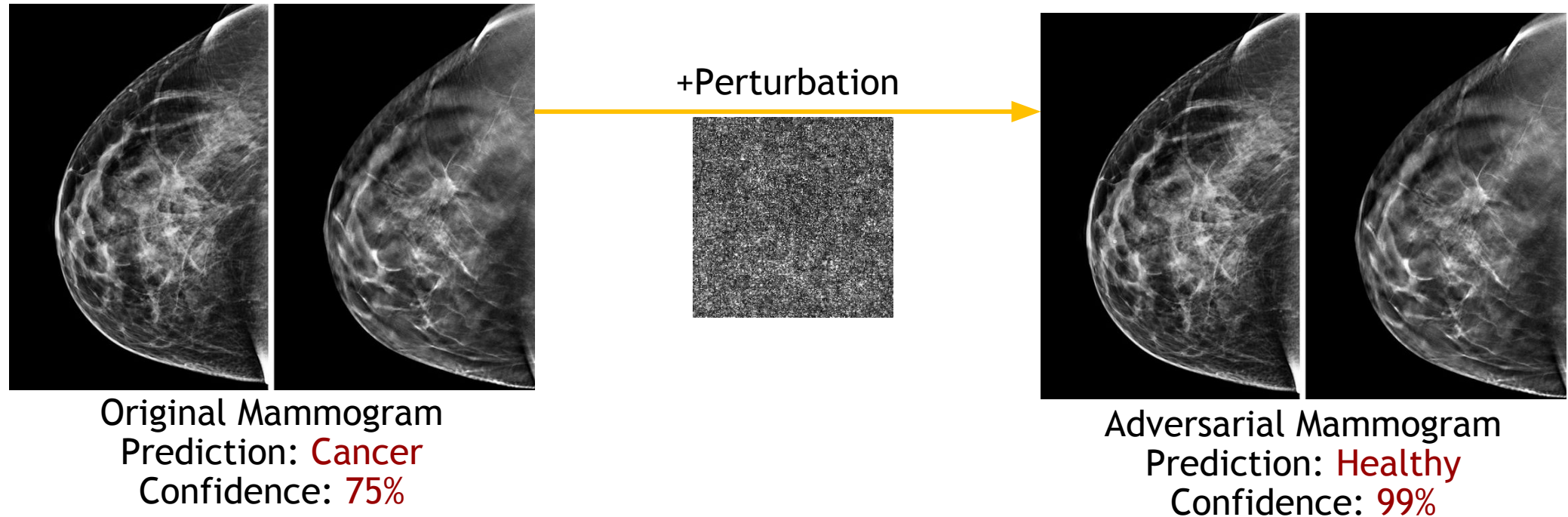
## Vehicle Detection





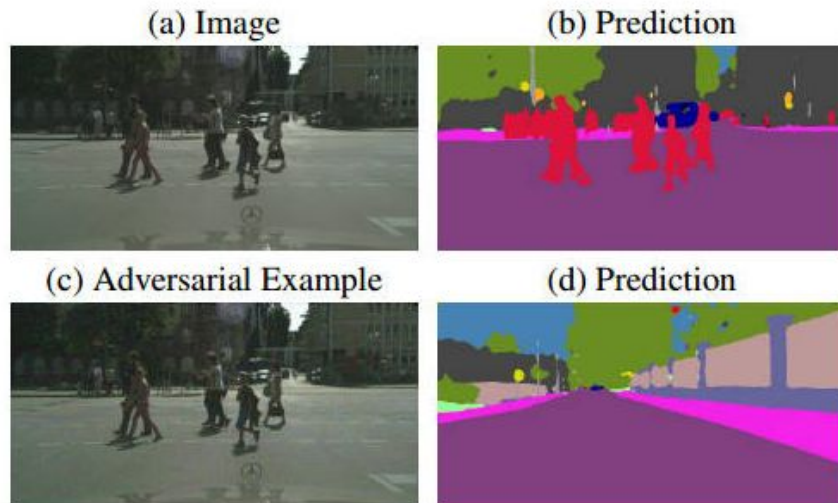
# Real World Consequences of Adversarial Examples

Adversarial examples with malicious intent reduce trust in automated systems (Healthcare).

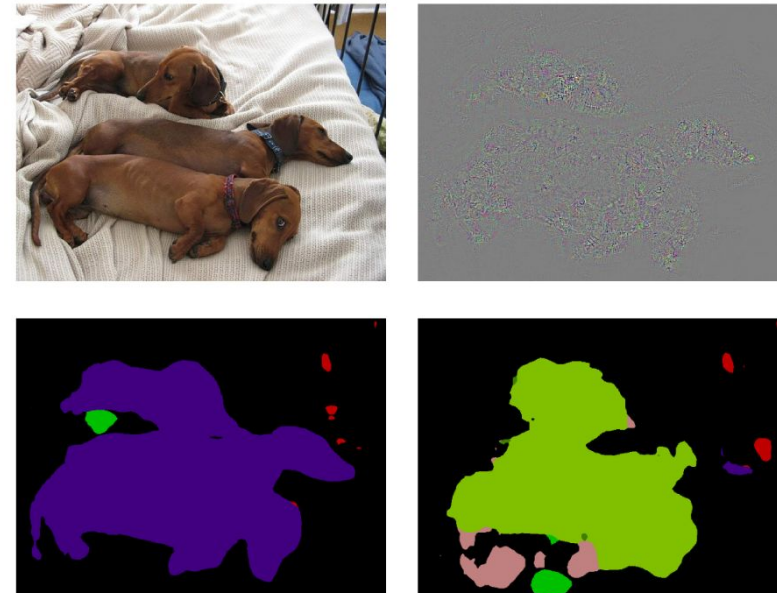


# Do Adversarial Examples Exist in Segmentation Models?

## Universal Perturbation (UP)



## Dense Adversary Generation (DAG)





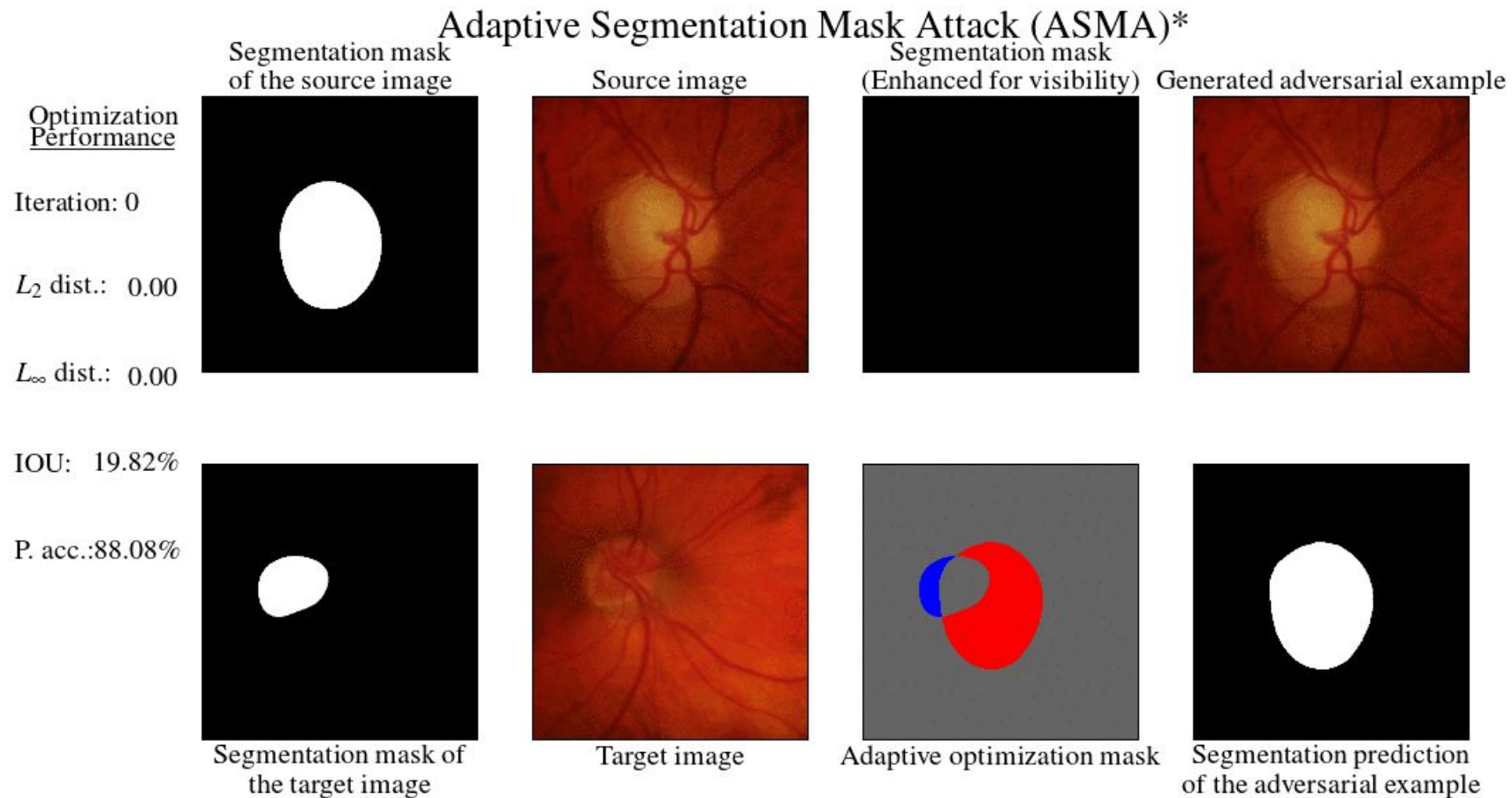
## Downsides of Existing Methods for Adversarial Example Generation for Segmentation

Segmentation models provide more information on the prediction than classification models. Thus, coming up with a defense for segmentation is much easier than classification (based on the shape of the prediction).

Existing attacks do not pose a threat because:

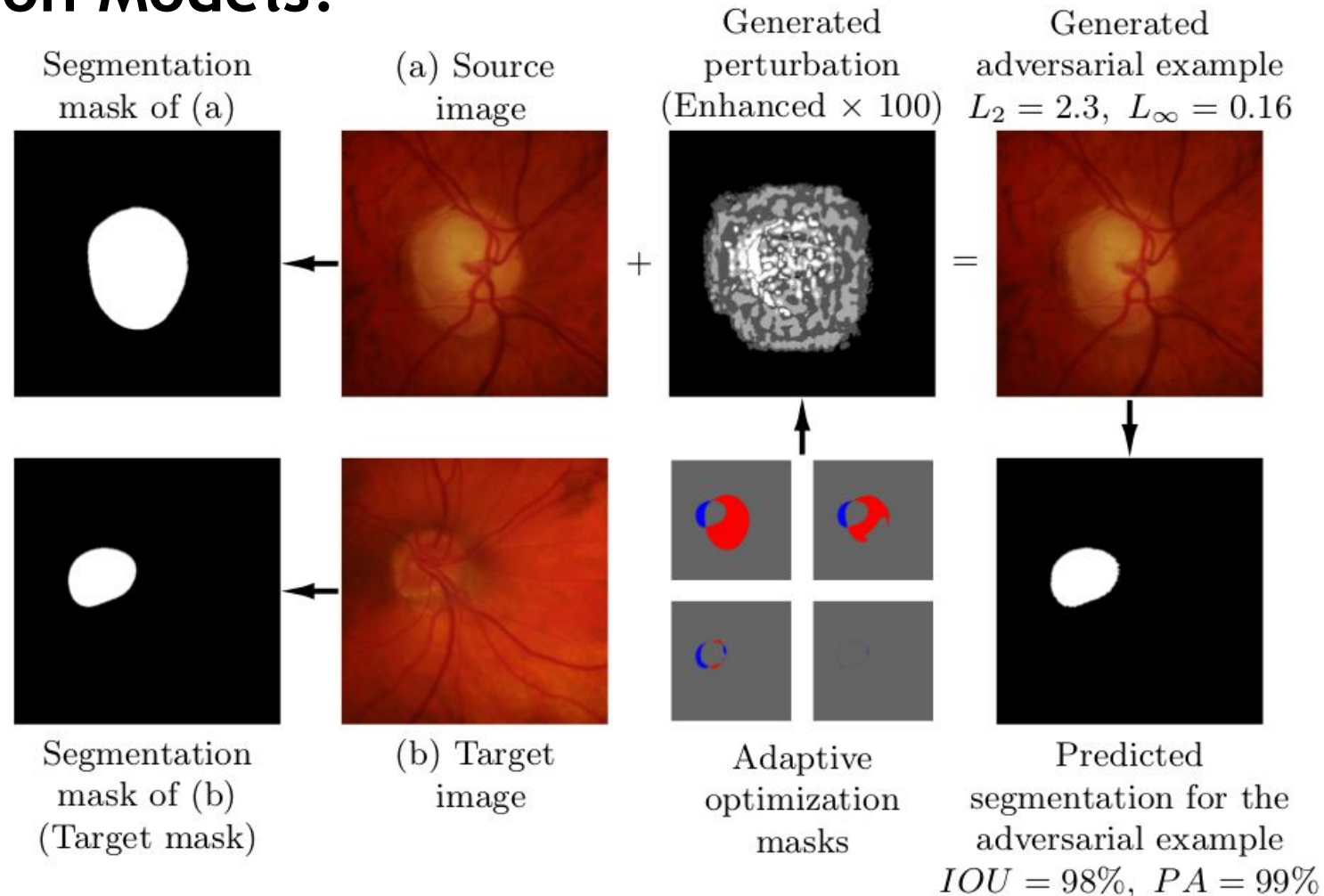
- DAG aims to misclassify all pixels, which leads to random-looking segmentations shapes.
- UP is proposed as a method to change segmentation prediction to a single target mask.

# So... Are Adversarial Examples are a threat to Segmentation Models?



\* Impact of Adversarial Examples on Deep Learning Models for Biomedical Segmentation, U. Ozbek et al. 22nd International Conference on Medical Image Computing and Computer Assisted Intervention, MICCAI-19

# So... Are Adversarial Examples are a threat to Segmentation Models?



## Quantitative Results of Our Approach

Optimization	Glaucoma Dataset				ISIC Skin Lesion Dataset			
	Modification		Accuracy		Modification		Accuracy	
	$L_2$	$L_\infty$	IoU	PA	$L_2$	$L_\infty$	IoU	PA
SSM	<b>4.60</b>	0.22	<b>47%</b>	94%	<b>11.76</b>	0.24	<b>43%</b>	88%
	$\pm 1.76$	$\pm 0.09$	$\pm 18\%$	$\pm 2\%$	$\pm 4.11$	0.05	$\pm 15\%$	$\pm 2\%$
ASM	2.82	0.17	94%	99%	4.11	0.16	89%	98%
	$\pm 1.29$	$\pm 0.09$	$\pm 7\%$	$\pm 1\%$	$\pm 2.23$	$\pm 0.10$	$\pm 9\%$	$\pm 1\%$
ASM + DPM (ASMA)	<b>2.47</b>	0.17	<b>97%</b>	99%	<b>3.88</b>	0.16	<b>89%</b>	98%
	$\pm 1.05$	$\pm 0.09$	$\pm 2\%$	$\pm 1\%$	$\pm 1.99$	$\pm 0.09$	$\pm 10\%$	$\pm 1\%$

**Thank you for listening!**

**Any questions?**

**Adversarial potato!**

