

Transferability in Machine Learning: from Phenomena to Black-Box Attacks using Adversarial Samples



भारतीय प्रौद्योगिकी
संस्थान जम्मू
INDIAN INSTITUTE OF
TECHNOLOGY JAMMU



Saurabh Gajbhiye
Harshit Shakya
Harsh Sharma
Rahul Kumar
Sachin Pal

21 April 2021

Contents

1. Introduction
2. Key terms
3. White box attack
4. Crafting Adversarial Attacks
5. Blackbox Attack and Transferability
6. Results
7. Conclusion

Introduction

Many machine learning models are vulnerable to
Adversarial Examples attack

We will briefly discuss how these black box attacks are crafted

We will see how the result of these attack on real world machine
learning models

Key Term

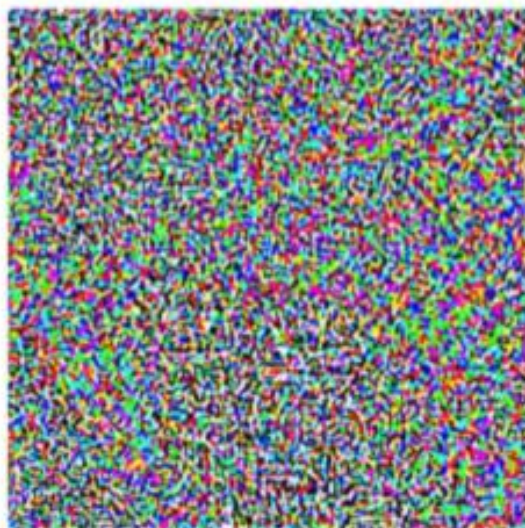
Adversarial Example

Inputs that are specially crafted to cause a machine learning model to produce an incorrect output



“panda”
57.7% confidence

$+ .007 \times$



Noise

$=$



“gibbon”
99.3 % confidence

ML application and consequence

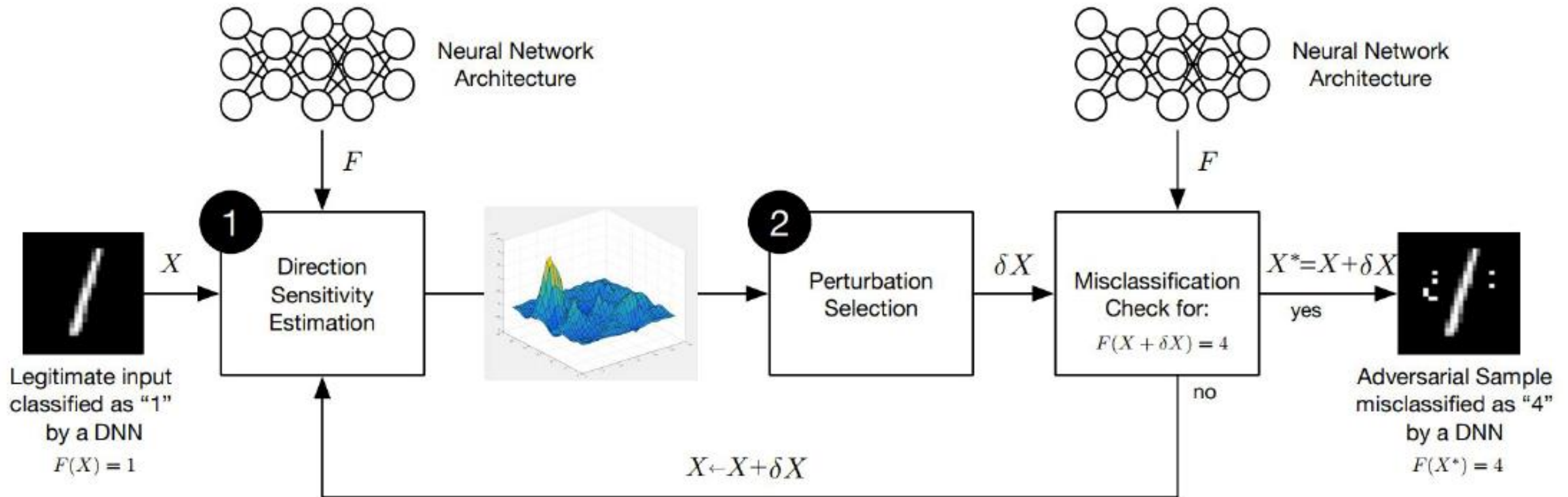
Autonomous Driving

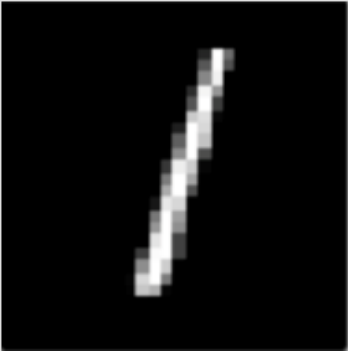
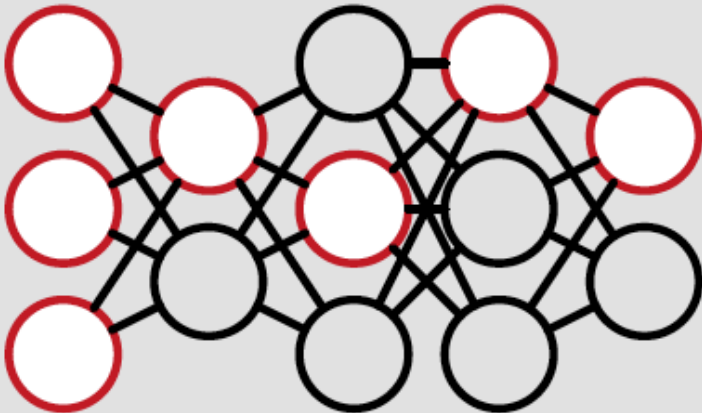

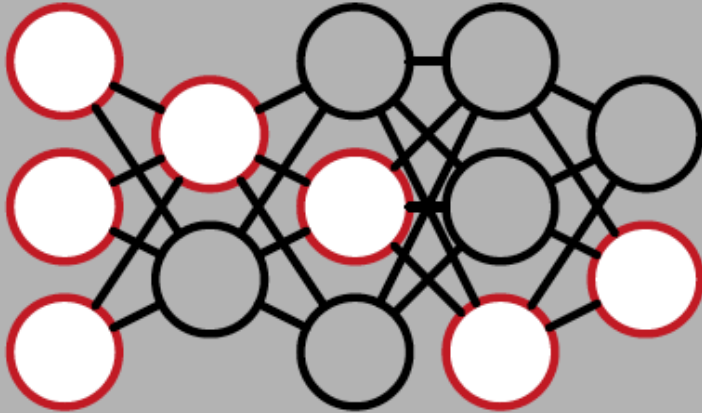
Fraud Detection in Finance

Malware Detection

Machine Learning as a Service Platform

White Box Attack



	Input	Model Activations	Output
Legitimate			1
Adversarial			4

Crafting Adversarial Attacks

DNN , LR and kNN use the **fast gradient sign method**

This is a white box attack

Attacker is aware about the :

Model Architecture

Training Data

Model Parameters

Fast gradient sign method

During training, the classifier uses a loss function to **minimize** model prediction errors

After training, **attacker** uses loss function to **maximize** model prediction error

1. Compute its gradient with respect to the input of the model
2. Take the sign of the gradient and multiply it by a threshold

Black box attack

In this attack the attacker is not aware about:

Model Architecture

Training Data

Model Parameters

He has access to model via API, when a input is passed to a model, a label (output) is returned.

Threat Model of Black box attack

Adversarial capabilities

~~Training data
Model architecture
Model parameters
Model scores~~



(limited) oracle
access: *labels*

Adversarial goal

Force a ML model remotely accessible through an API to misclassify

Example



Challenges in Black Box attack

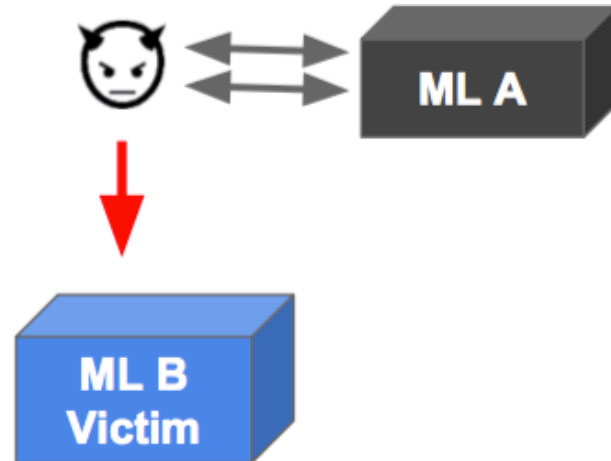
Alleviate lack of knowledge
about model

Alleviate lack of
training data

Challenge 1: Use Adversarial example transferability

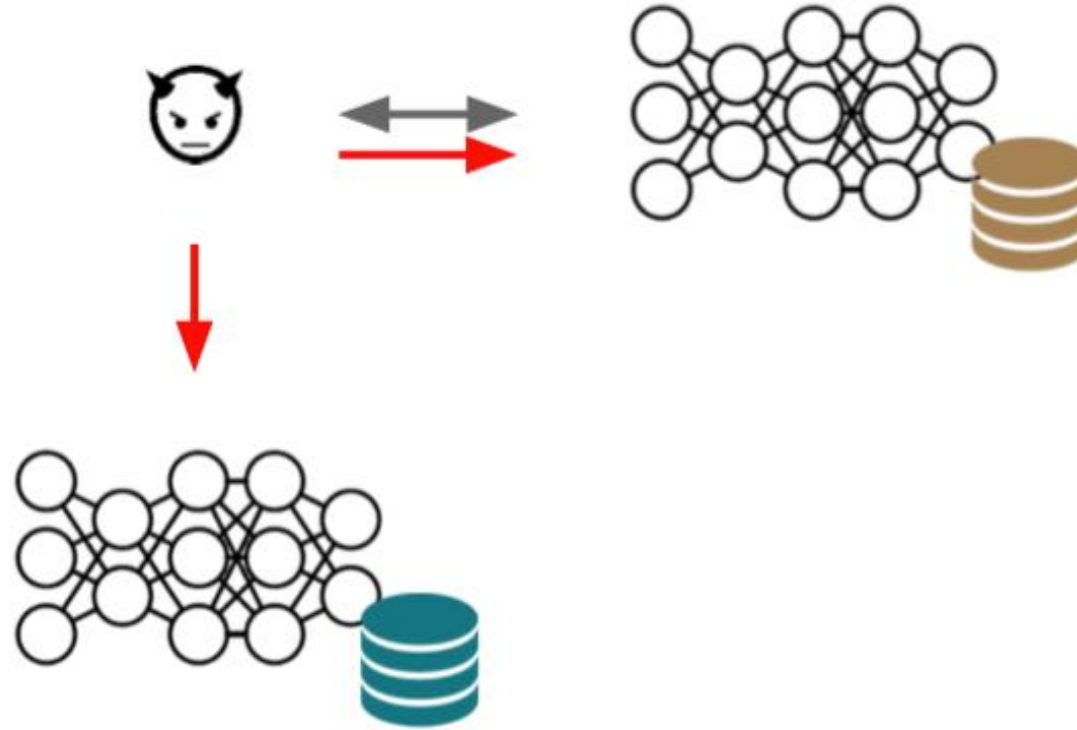
Adversarial examples have a **transferability** property:

samples crafted to mislead a model A are likely to mislead a model B



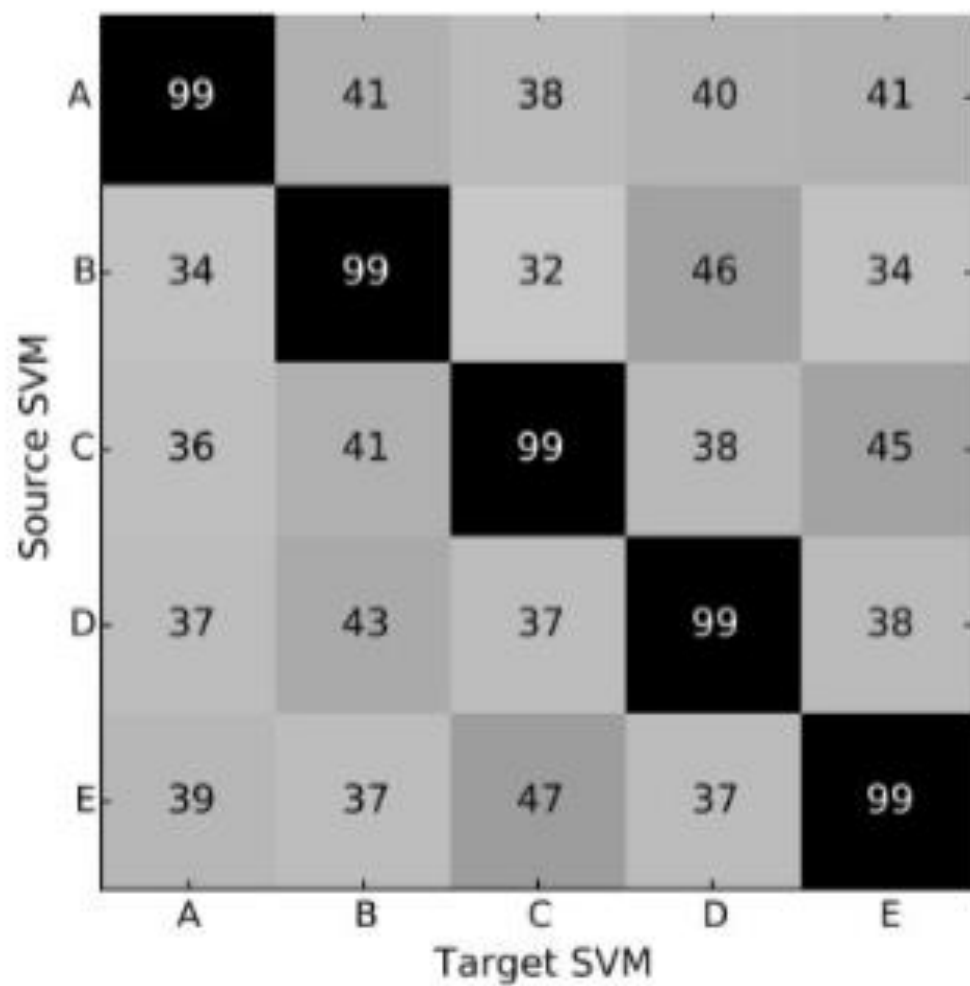
Intra-technique Transferability

samples crafted to mislead a model A are likely to mislead a model B

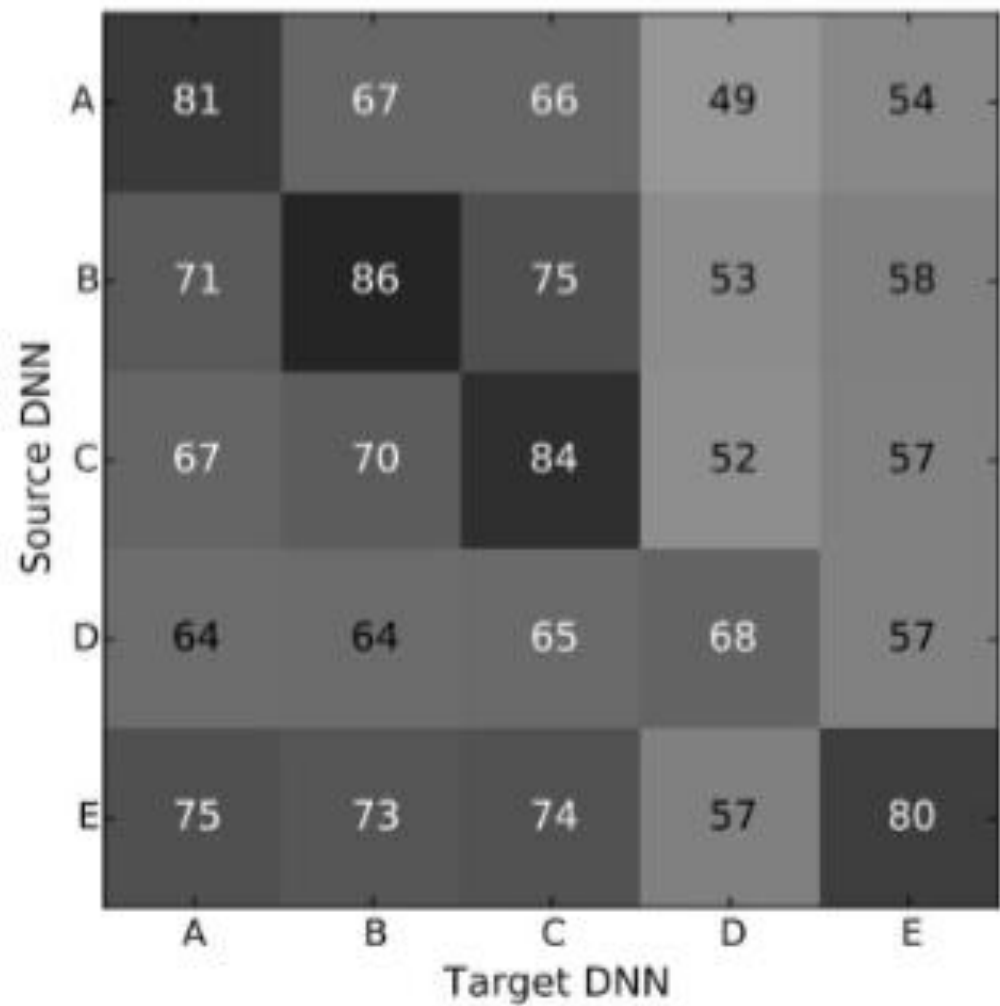


Source LR	A	B	C	D	E
	98	95	95	95	95
	95	98	95	95	94
	94	94	98	95	95
	94	95	95	98	95
	95	95	95	95	98
Target LR					

Strong



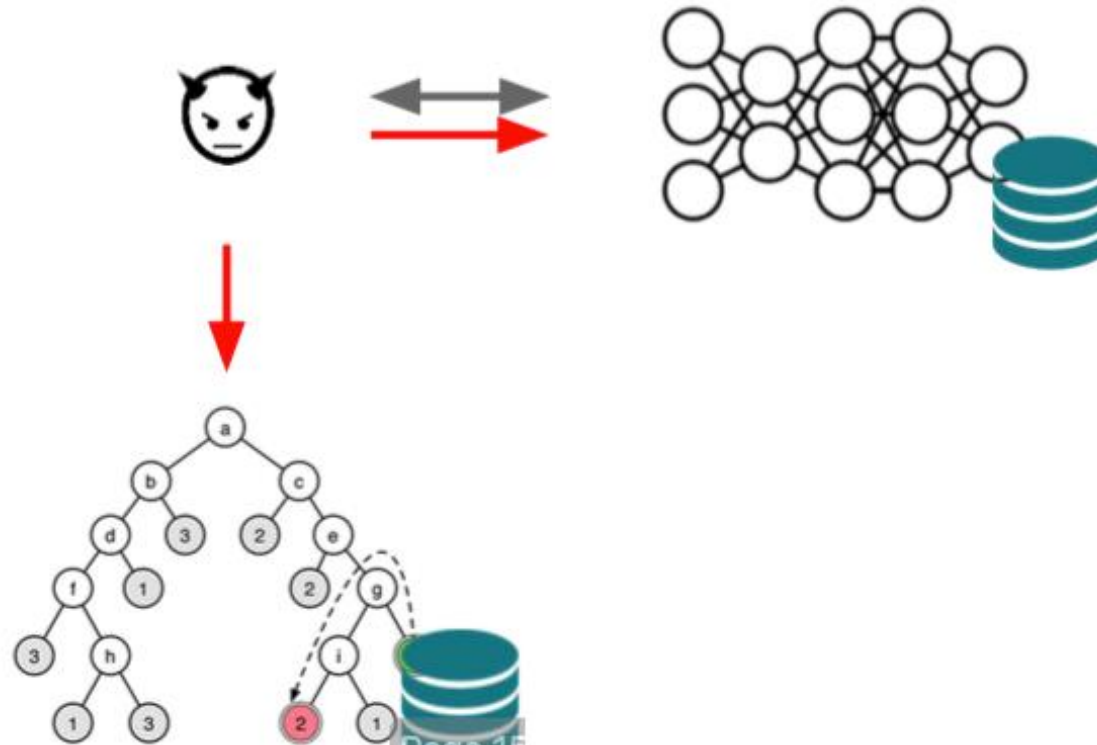
Weak

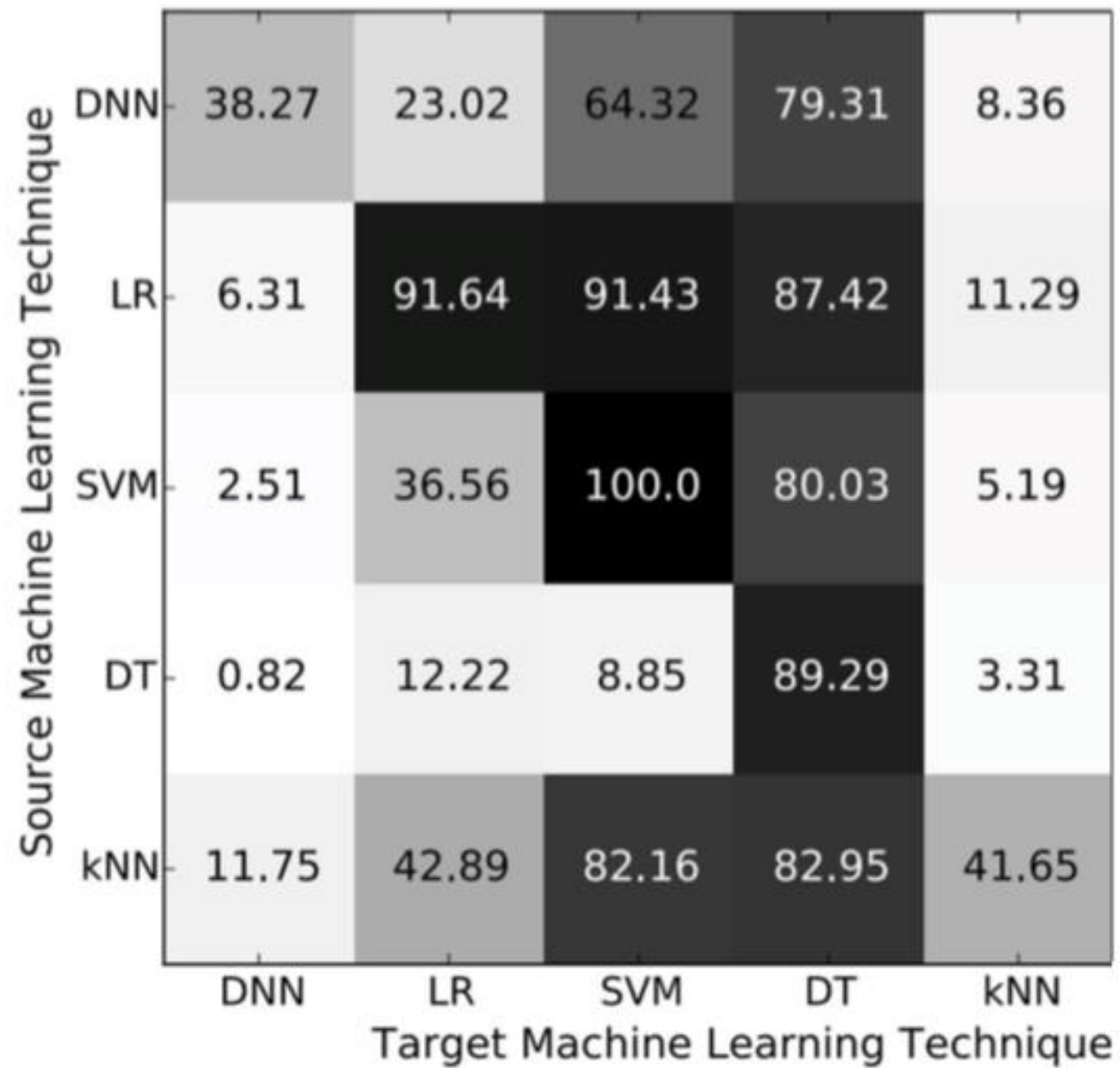


Intermediate

Cross-technique Transferability

samples crafted to mislead a model A are likely to mislead a model B





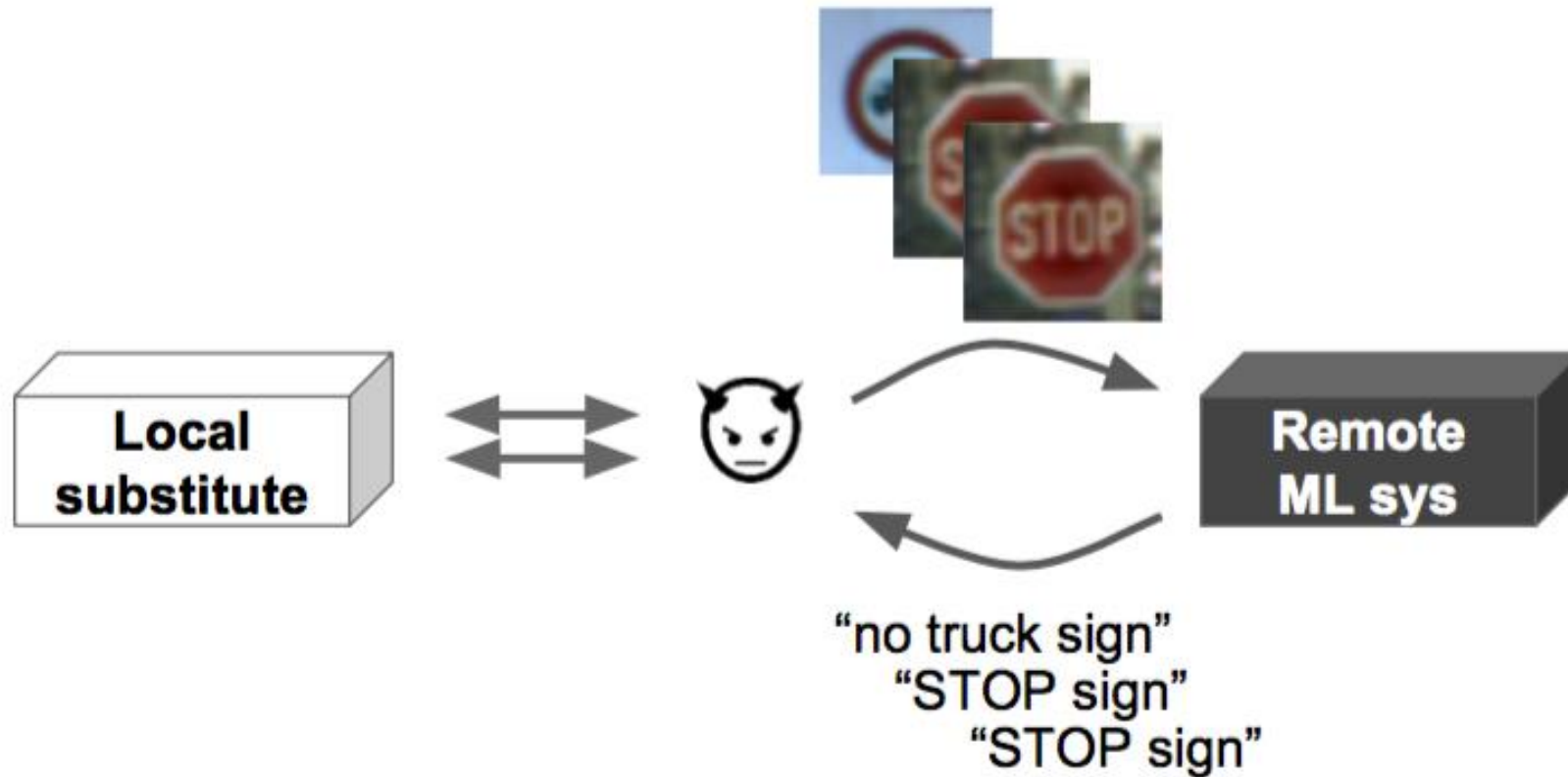
Alleviate lack of knowledge
about model

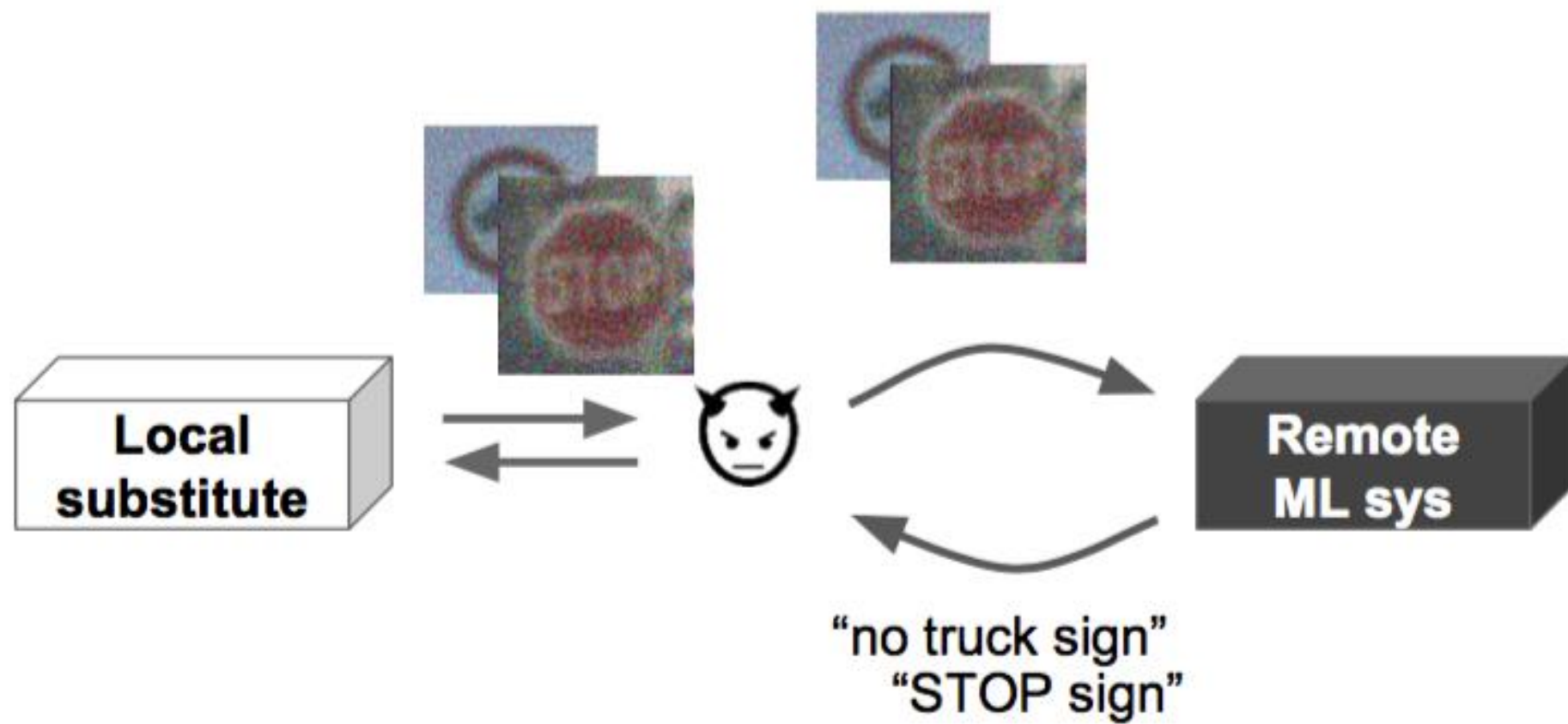


Adversarial example
transferability from a
substitute model to
target model

Alleviate lack of
training data

Challenge 2: Training Model using Synthetic data generation







Alleviate lack of knowledge
about model



Adversarial example
transferability from a
substitute model to
target model




+

Alleviate lack of
training data



Synthetic data
generation

Results

Remote Platform	ML technique	Number of queries	Adversarial examples misclassified (after querying)
 MetaMind	Deep Learning	6,400	84.24%
 amazon web services™	Logistic Regression	800	96.19%
 Google Cloud Platform	Unknown	2,000	97.72%

Conclusion

We saw phenomenon of adversarial sample transferability across the machine learning models.

We demonstrated how to create black box attack that could be used to target online ML models trained and hosted by Amazon and Google, without any knowledge of the model design or parameters

Thank you!