

Using Adversarial Defense Methods to Improve the Performance of Deep-Neural-Network-Controlled Automatic Driving Systems

Mike Camara

Automatic Driving Systems



Automatic Driving Systems

Safer roads



Automatic Driving Systems

Safer roads



Freeing up time



Automatic Driving Systems

Safer roads



Freeing up time



Mobility for all



Automatic Driving Systems

Modular pipeline



End-to-end pipeline



Automatic Driving Systems

Modular pipeline

Expensive sensors



Cameras
LiDAR
Radar
HD Maps
GNSS

End-to-end pipeline



Automatic Driving Systems

Modular pipeline

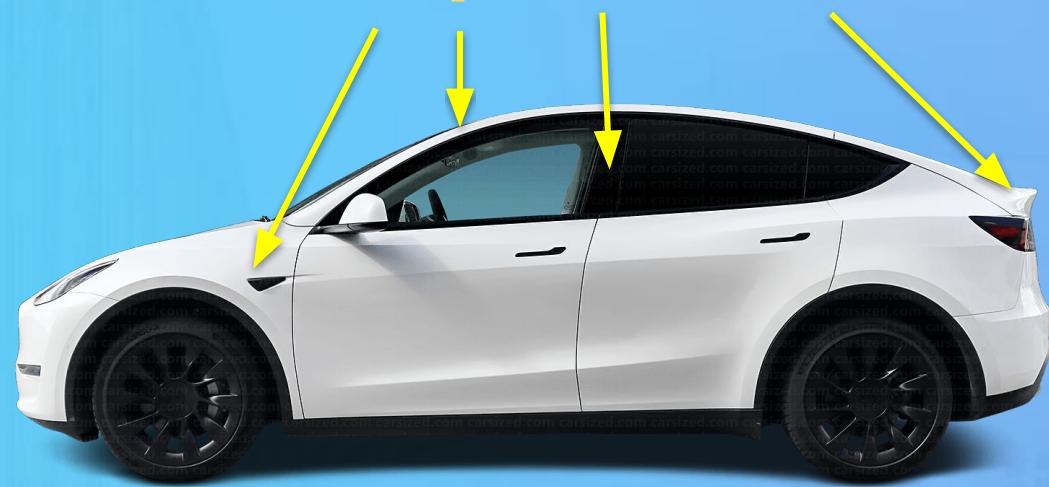
Expensive sensors



Cameras
LiDAR
Radar
HD Maps
GNSS

End-to-end pipeline

Computer Vision



Cameras
Data

Automatic Driving Systems

Modular pipeline

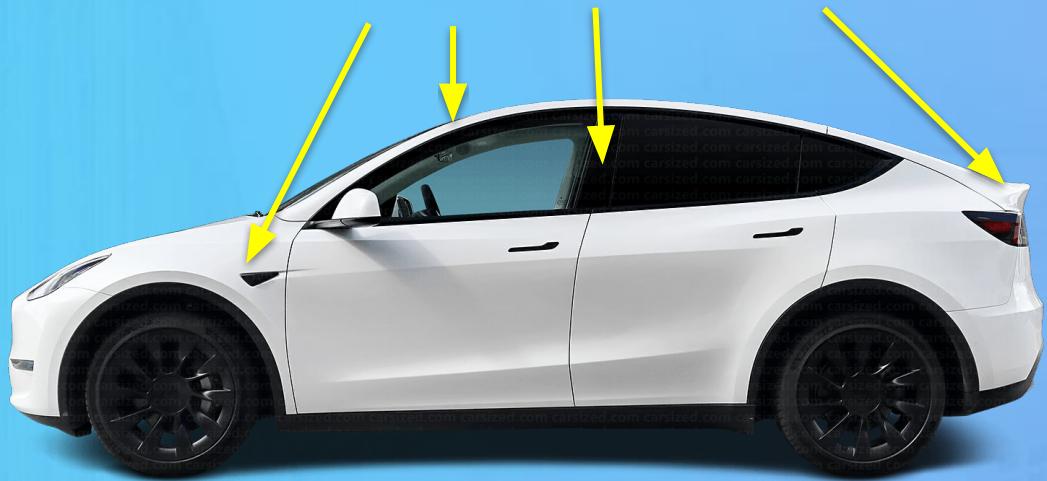
Expensive sensors



Cameras
LiDAR
Radar
HD Maps
GNSS

End-to-end pipeline

Computer Vision



Cameras

a lot of Data

End-to-end Pipeline

Collect

Labelled images



Dataset

End-to-end Pipeline

Collect

Labelled images

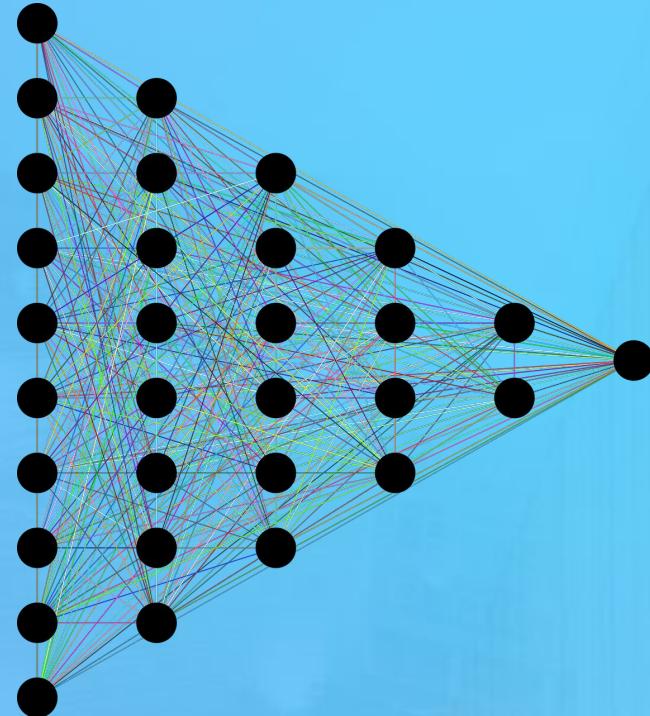


Dataset

**TAL
TECH**

Train

Deep learning model



Convolutional Neural Network

CNN

End-to-end Pipeline

Collect

Labelled images

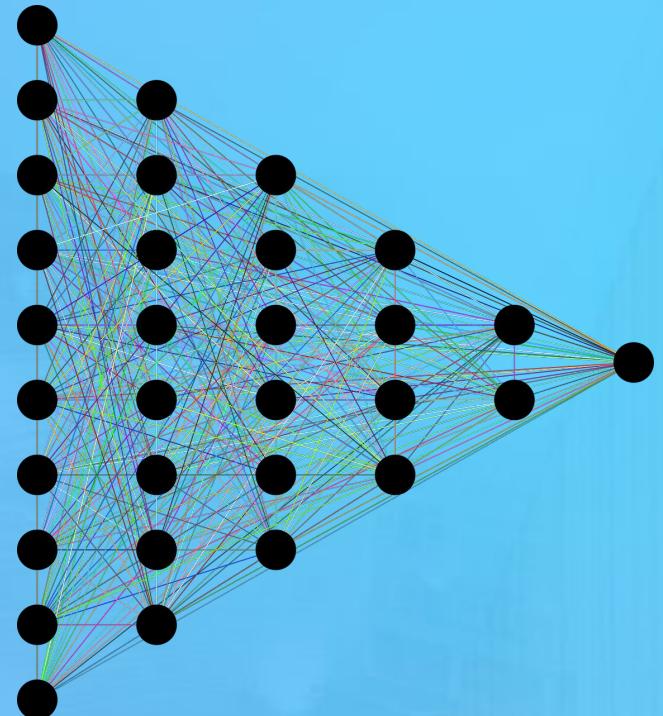


Dataset

**TAL
TECH**

Train

Deep learning model



CNN

Output

Model



autopilot.h5

End-to-end Pipeline

Collect

Labelled images

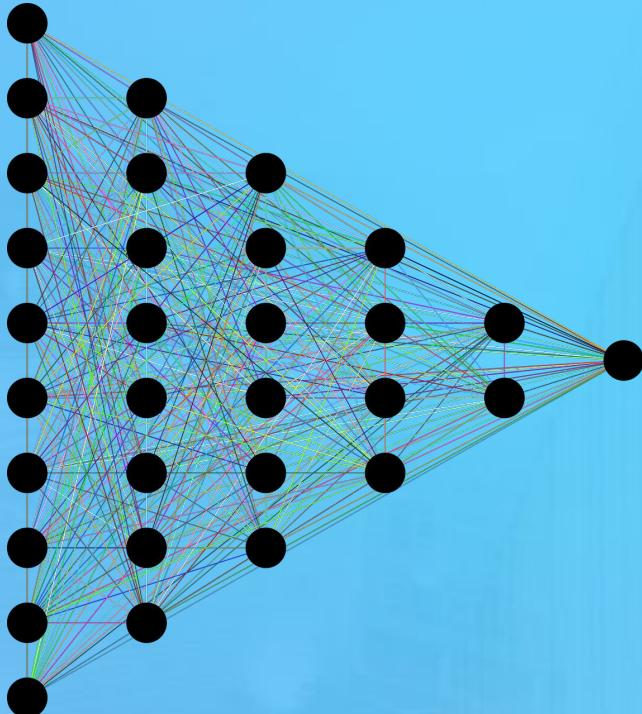


Dataset

TAL
TECH

Train

Deep learning model



Convolutional Neural Network

CNN

Output

Model



autopilot.h5

Deploy

Model in production



ADS
Steering
Acceleration
Brake

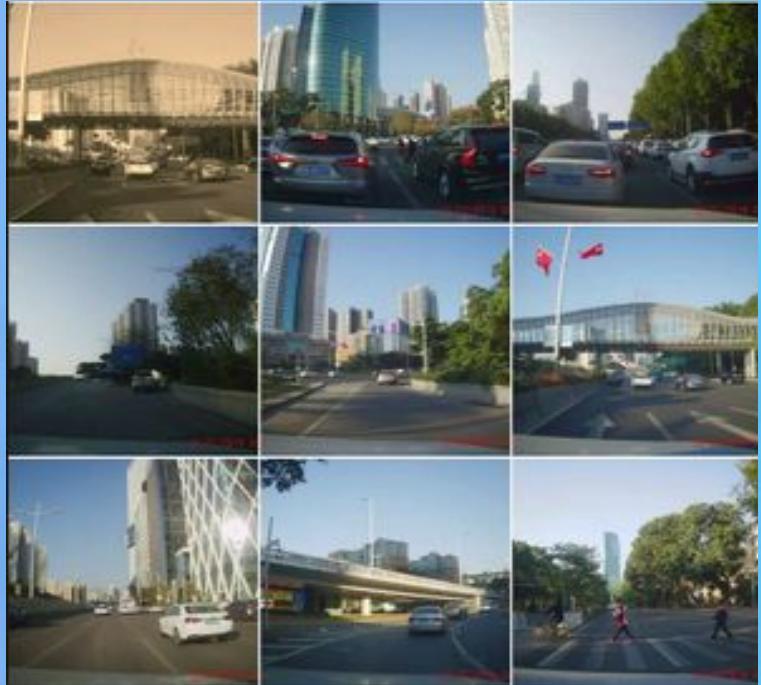


End-to-end Pipeline



Large dataset

End-to-end Pipeline



End-to-end Pipeline



Collect the dataset in
different lighting levels



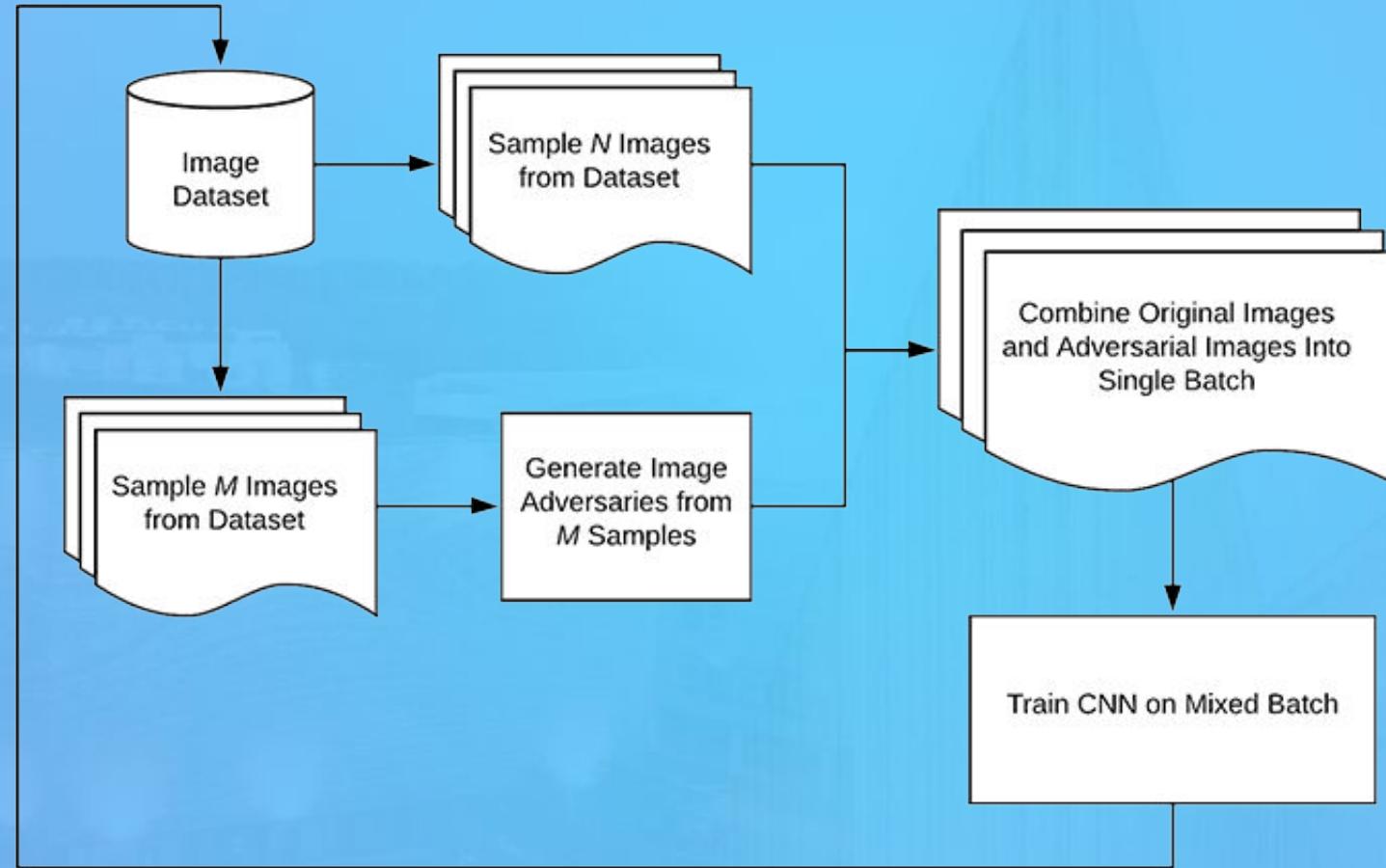
End-to-end Pipeline



Collect the dataset in
different lighting levels

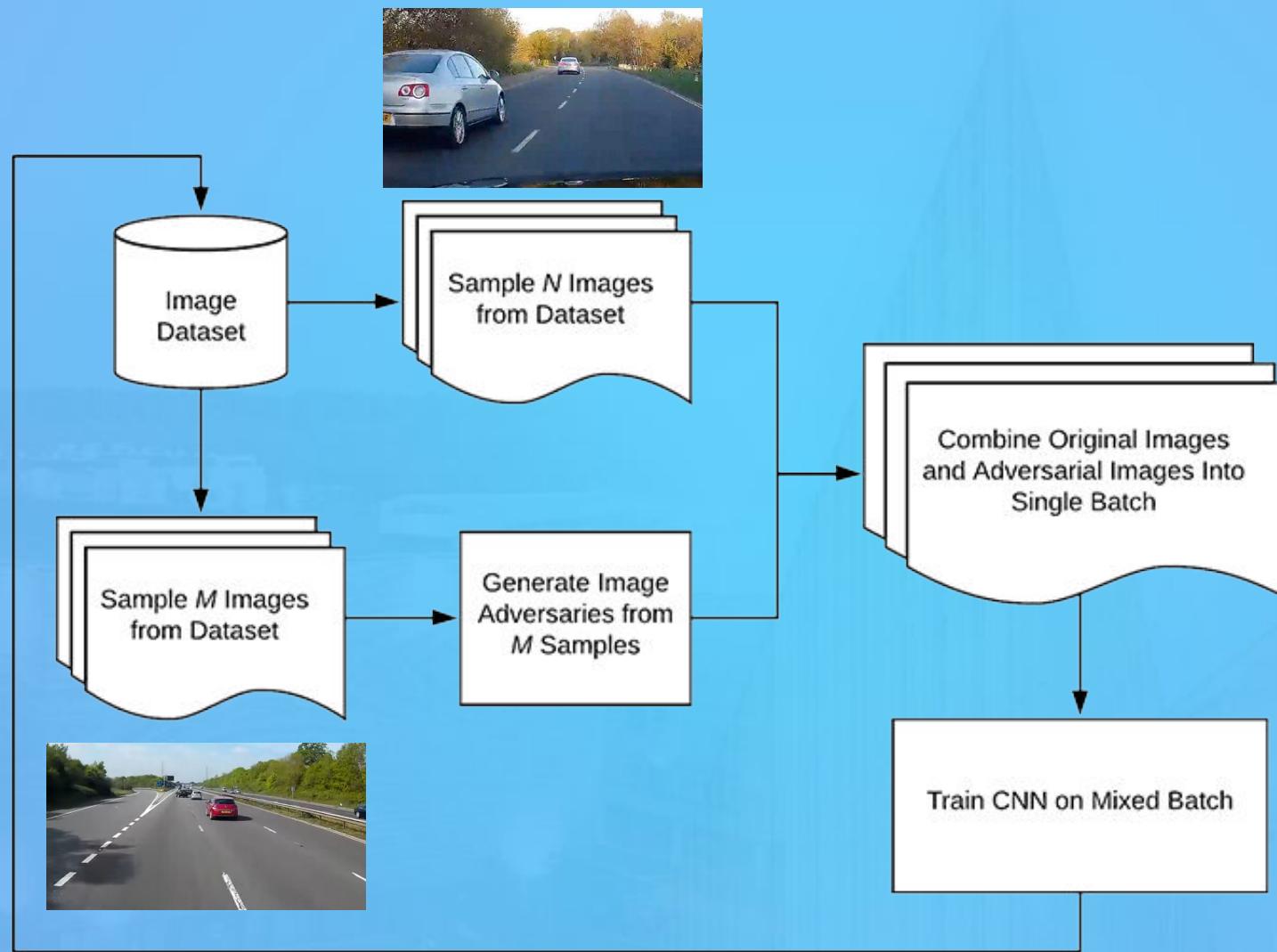


CNN's Adversarial Defense Training Method



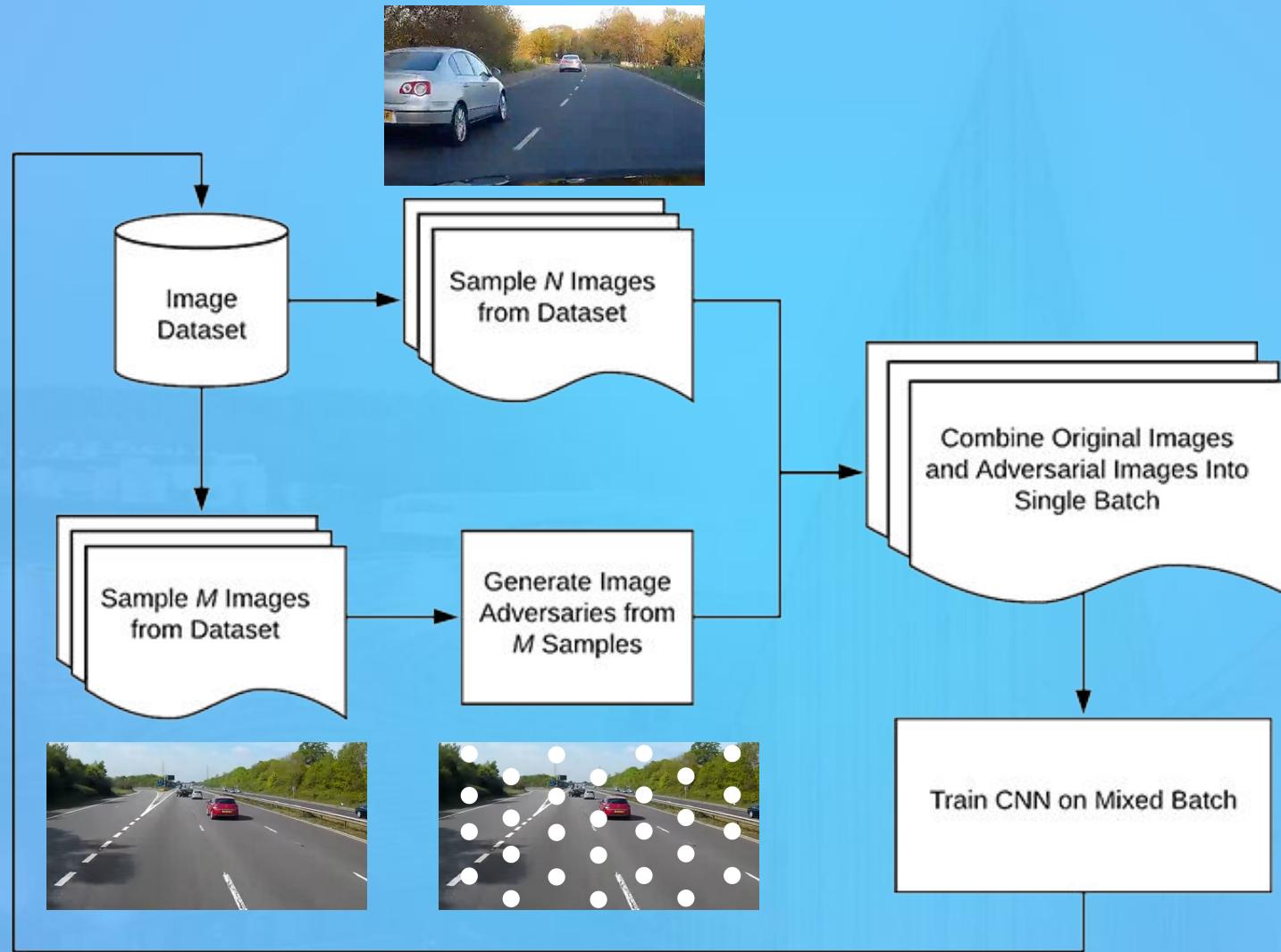
(Rosebrock, 2021)

CNN's Adversarial Defense Training Method



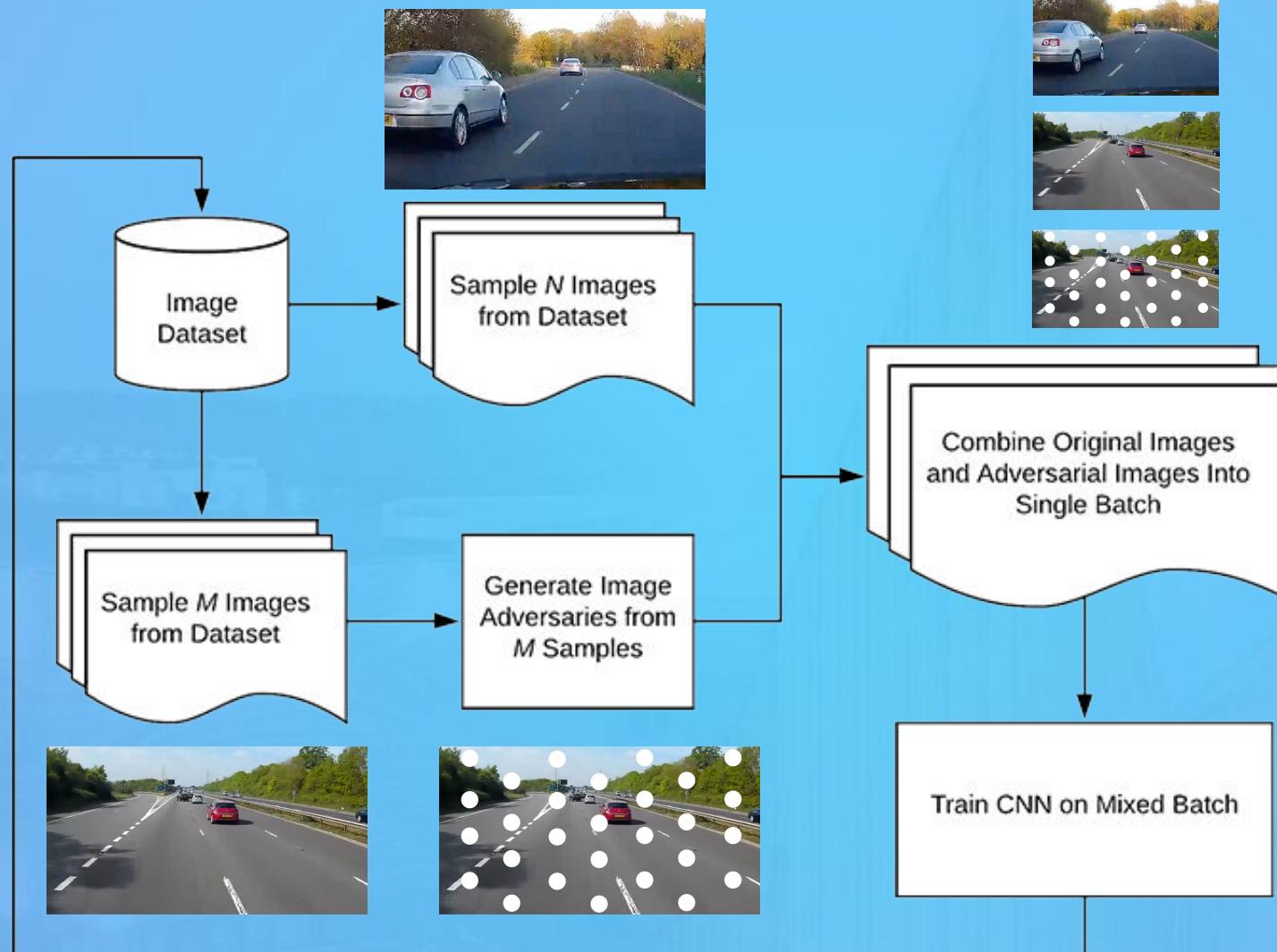
(Rosebrock, 2021)

CNN's Adversarial Defense Training Method



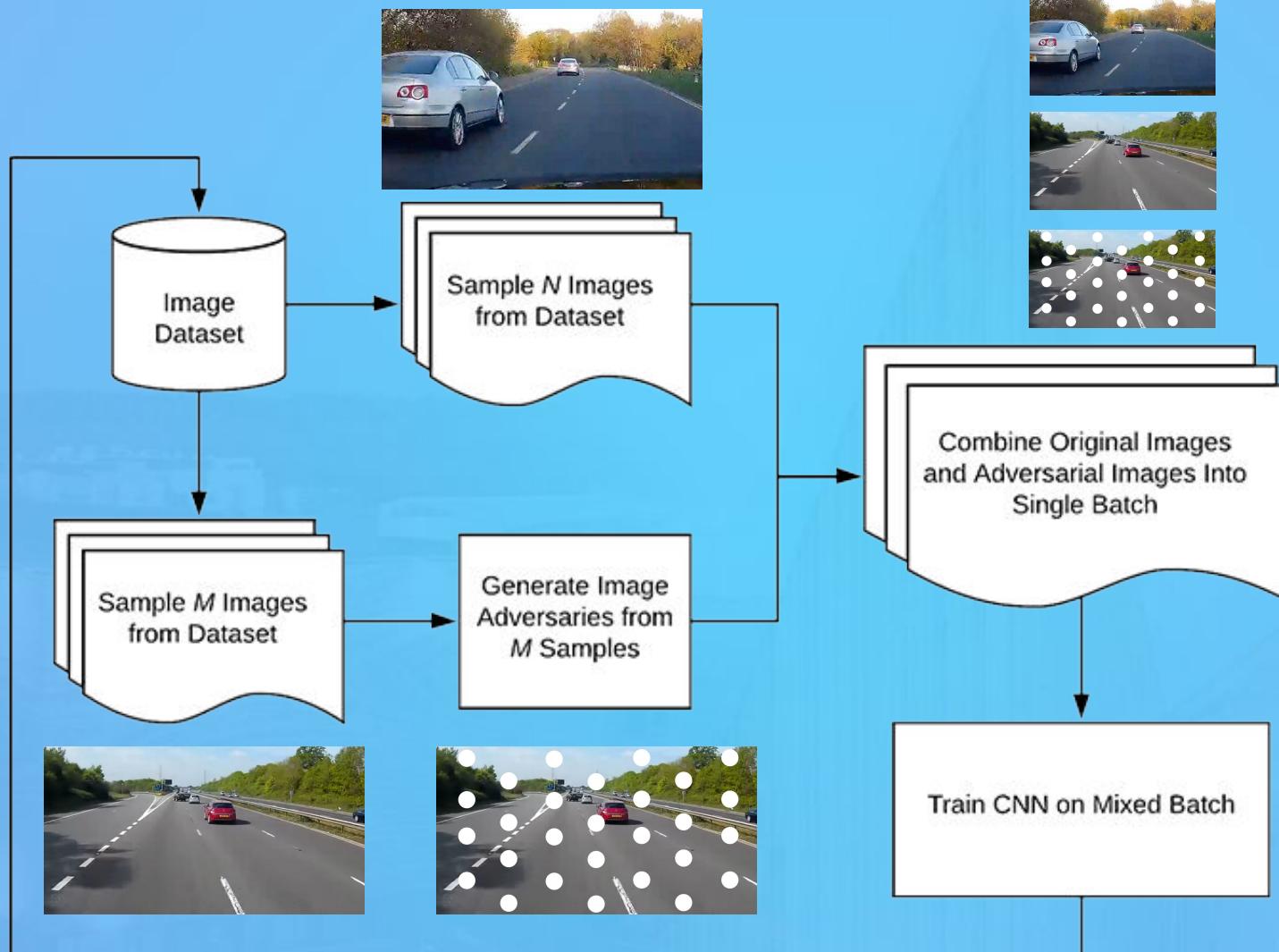
(Rosebrock, 2021)

CNN's Adversarial Defense Training Method

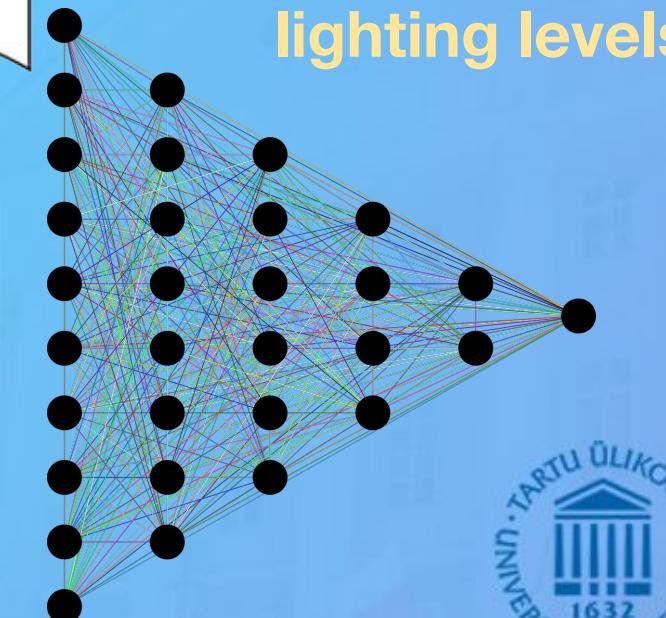


(Rosebrock, 2021)

CNN's Adversarial Defense Training Method



Model
better at
generalizing
to **unseen**
lighting levels



(Rosebrock, 2021)

Goal

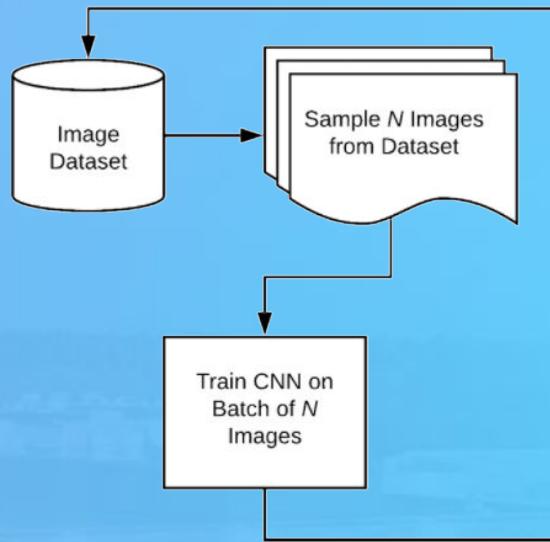
To **improve the performance** of CNN models
in unseen **lighting conditions**

Research Question

Can adversarial defense training methods
improve the **neural network** generalization skills
to unseen **lighting conditions**?

Hypothesis

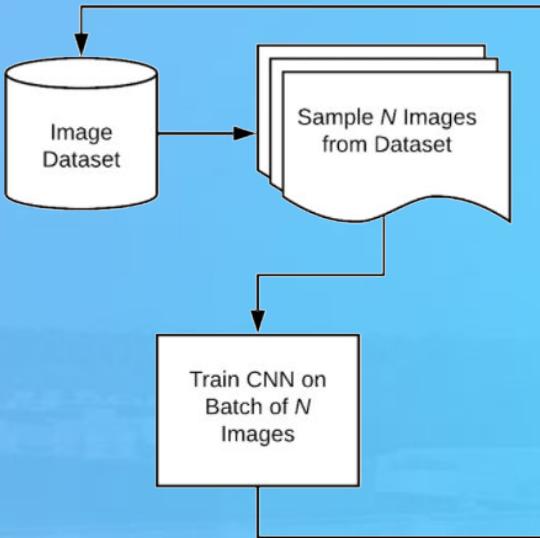
Train model



Conventional

Hypothesis

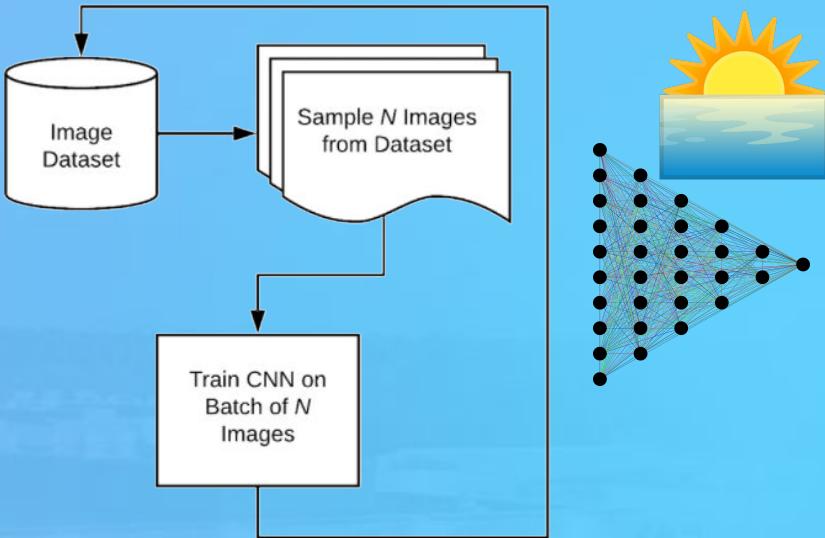
Train model



Conventional

Hypothesis

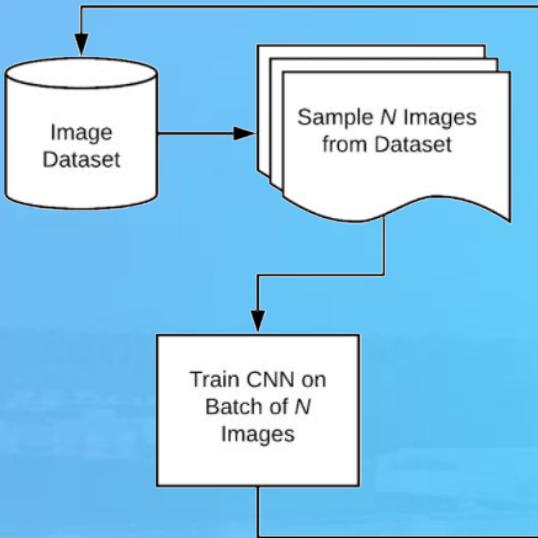
Train model



Conventional

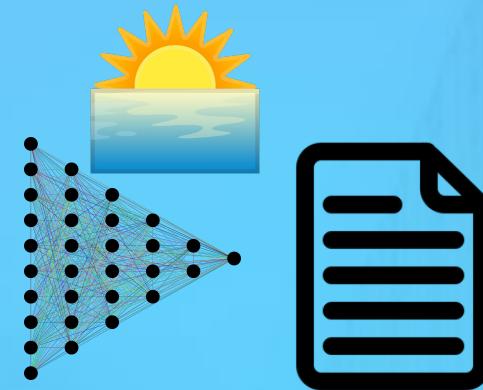
Hypothesis

Train model



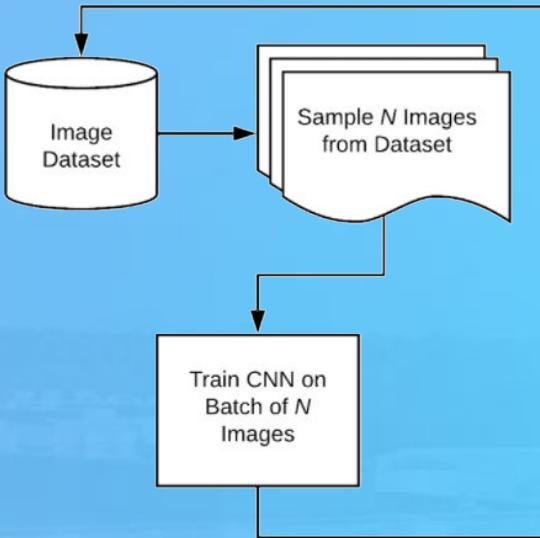
Conventional

Output model



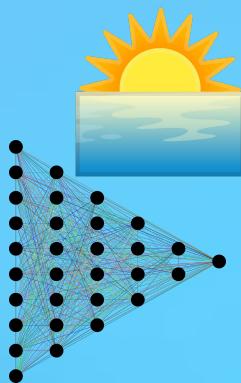
Hypothesis

Train model



Conventional

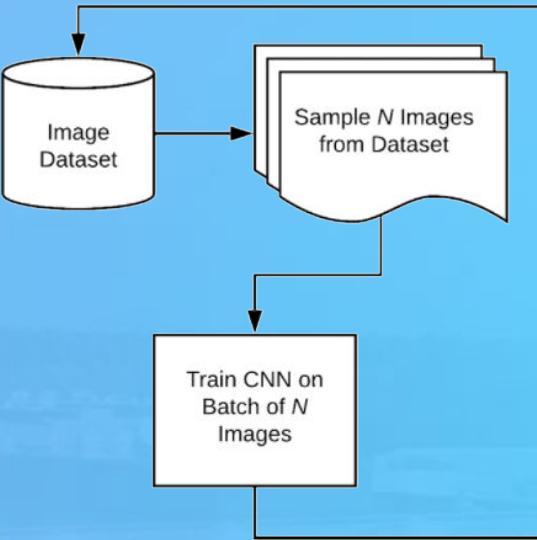
Output model



Deploy model in production

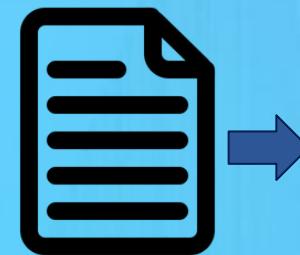
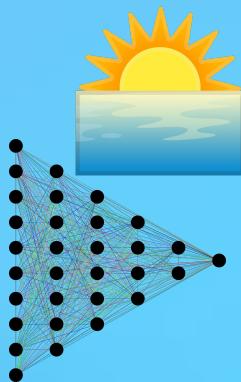
Hypothesis

Train model



Conventional

Output model

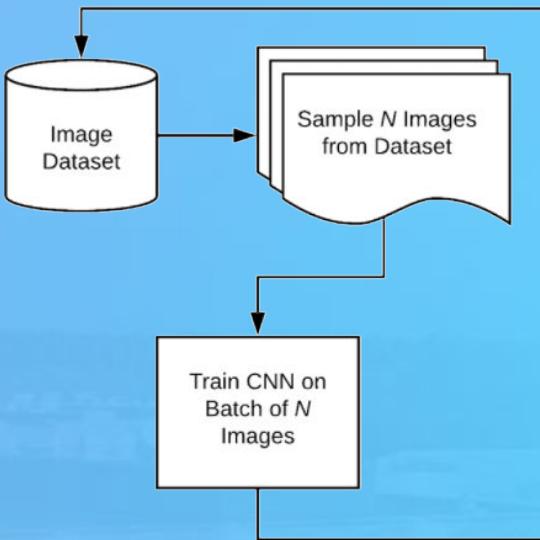


Deploy model in production



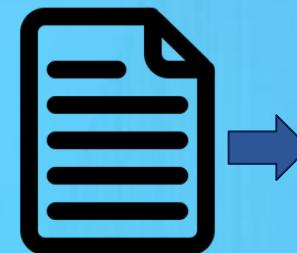
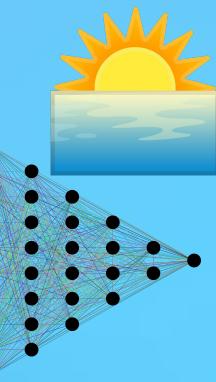
Hypothesis

Train model

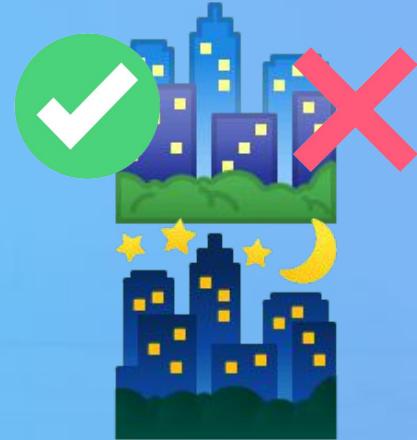


Conventional

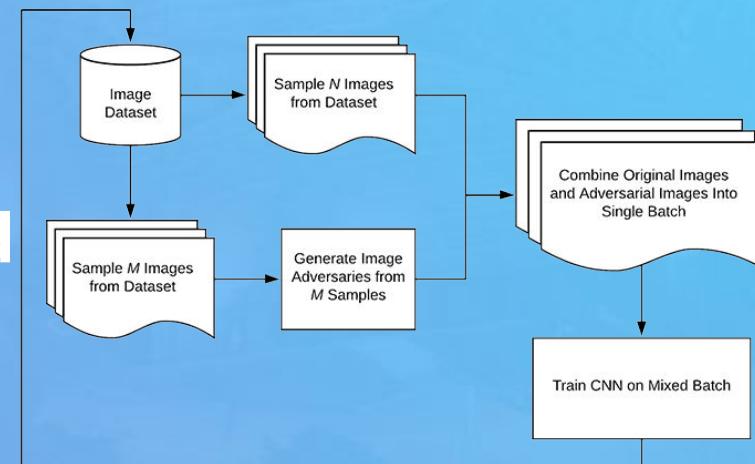
Output model



Deploy model in production

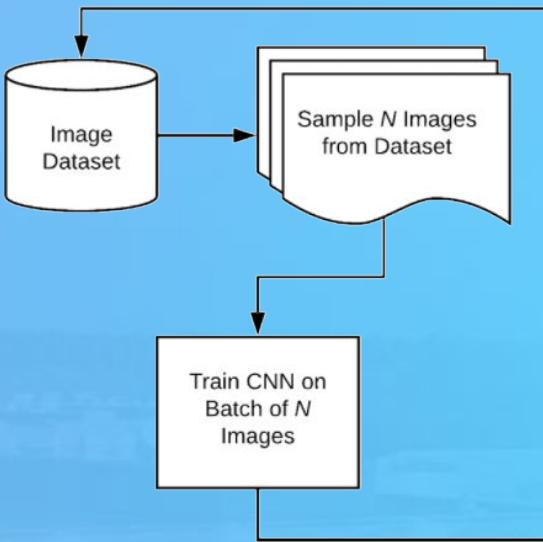


Proposed

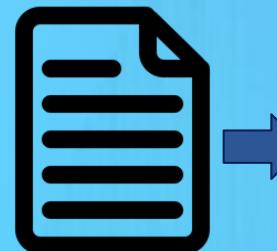


Hypothesis

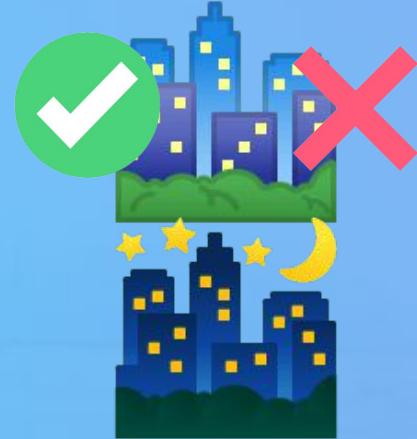
Train model



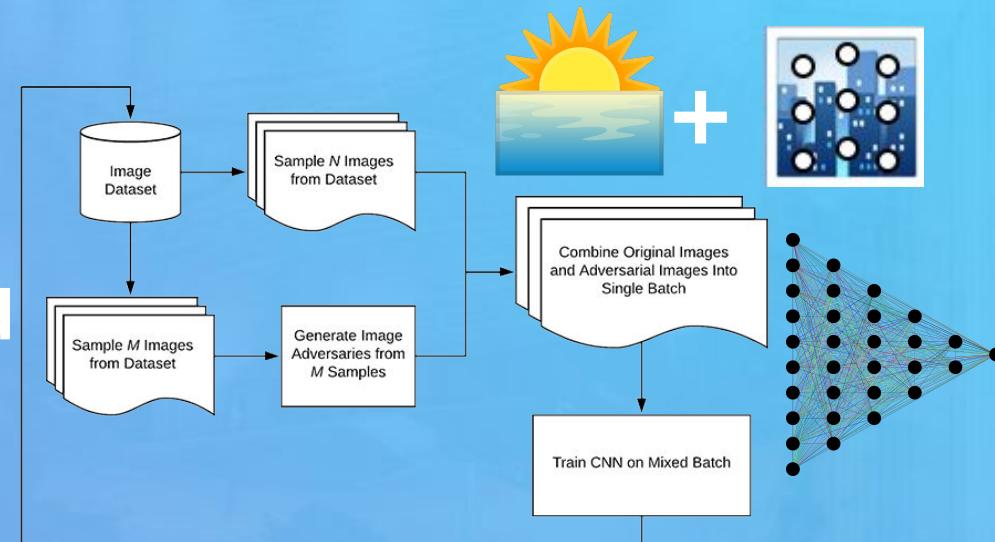
Output model



Deploy model in production



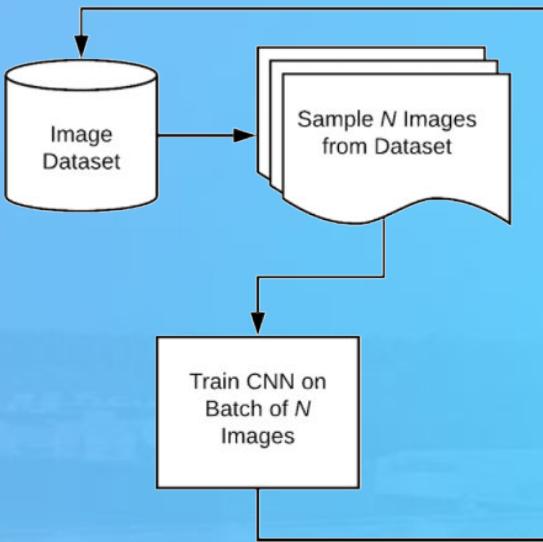
Conventional



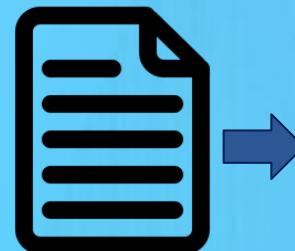
Proposed

Hypothesis

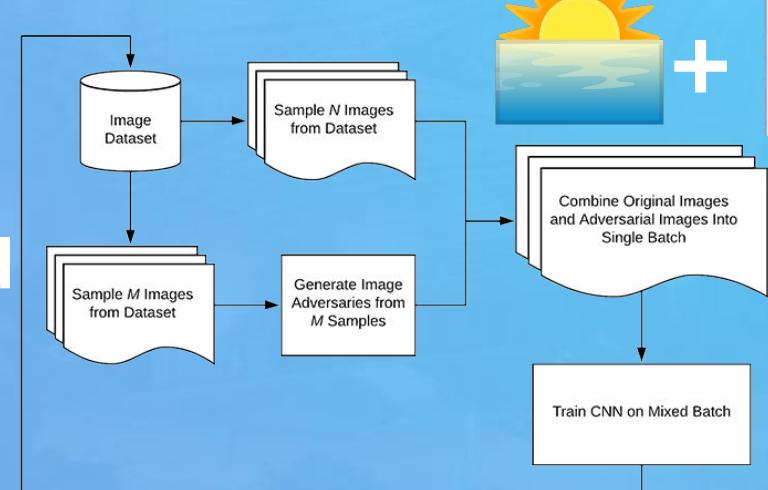
Train model



Output model



Conventional



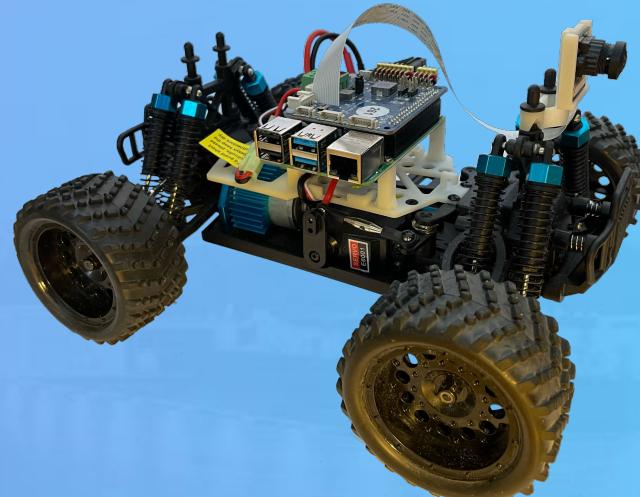
Proposed

Deploy model in production



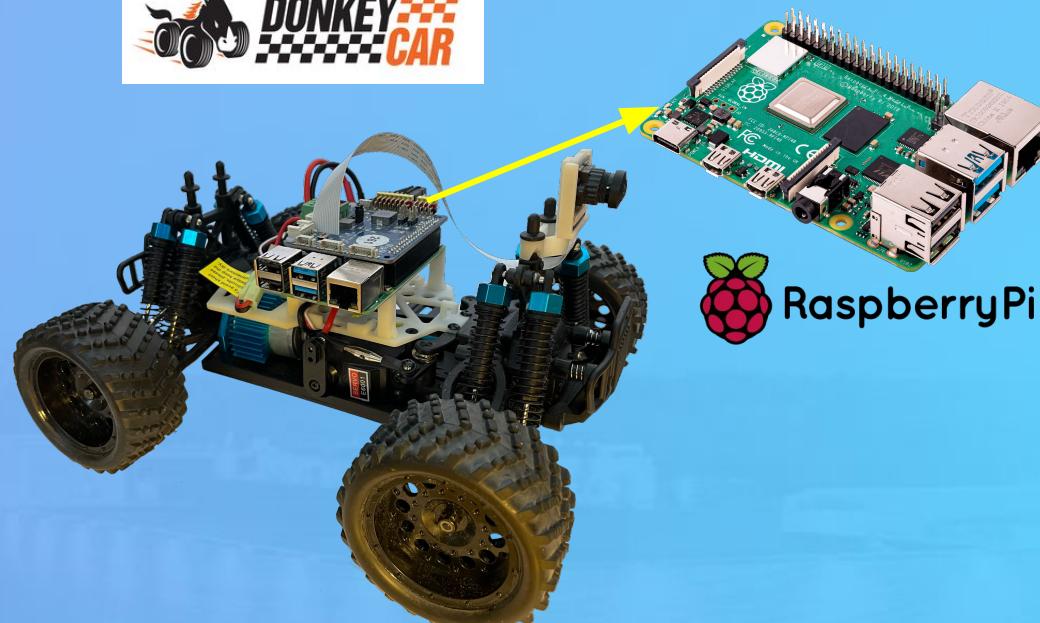
Method

Training CNN models



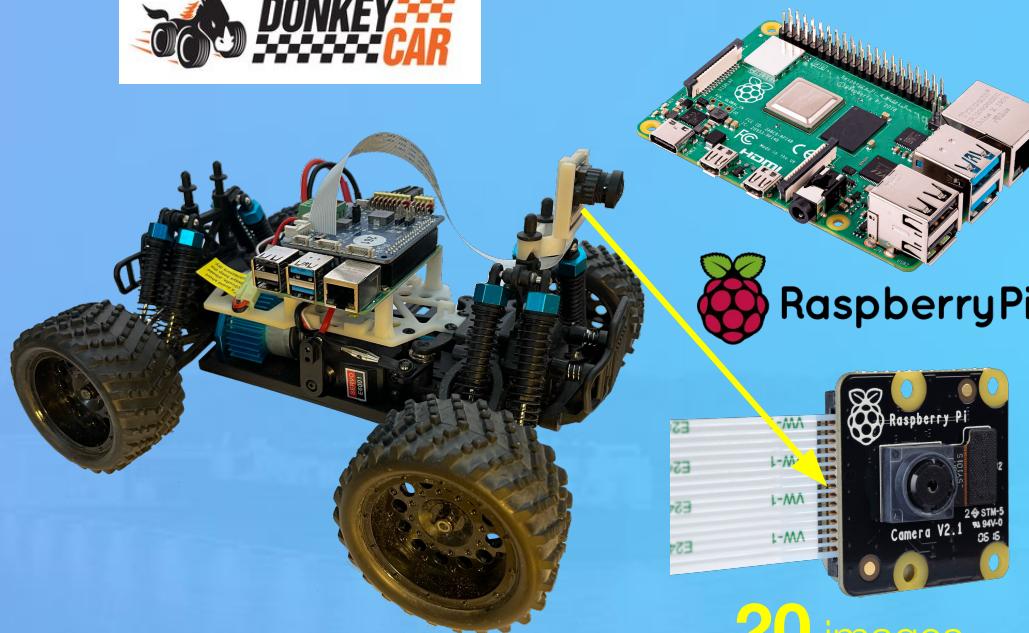
Method

Training CNN models



Method

Training CNN models



20 images
per second

Method



RaspberryPi



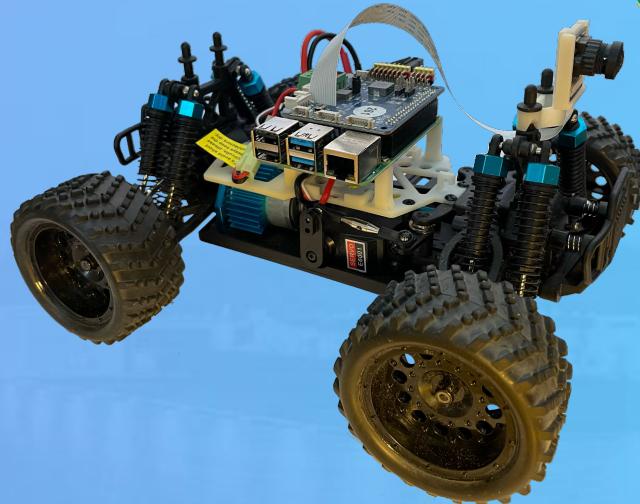
20 images
per second

Training CNN models

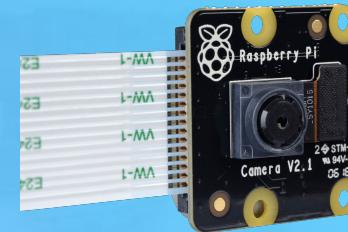
ONLY in lower-lights conditions



Method



RaspberryPi



20 images
per second

Training CNN models

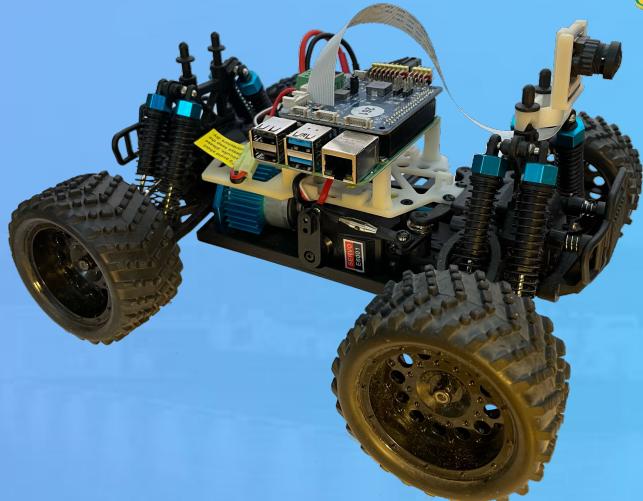
ONLY in lower-lights conditions



Imitation Learning

TAL
TECH

Method



RaspberryPi

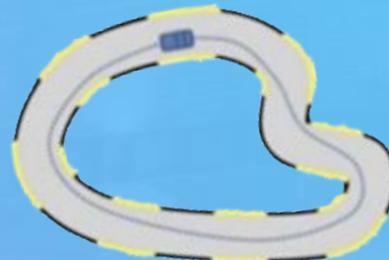


20 images
per second

Expert Demonstration

Imitation Learning

TAL
TECH

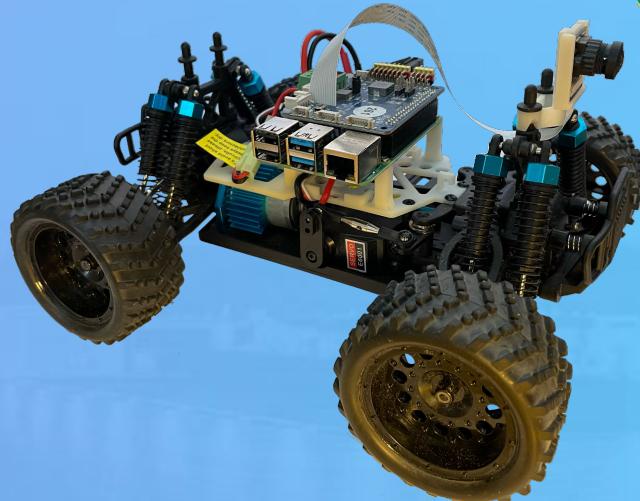


Training CNN models

ONLY in lower-lights conditions



Method



RaspberryPi

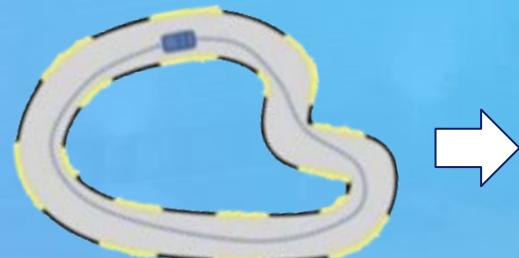


20 images
per second

Expert Demonstration

Imitation Learning

TAL
TECH



Training CNN models

ONLY in lower-lights conditions

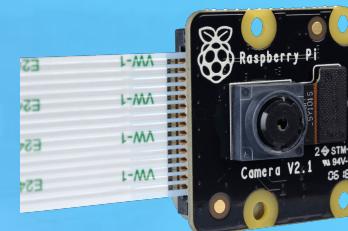


Dataset Cleaning

Method



RaspberryPi



20 images
per second

Expert Demonstration

Training CNN models

ONLY in lower-lights conditions



Imitation Learning

TAL
TECH

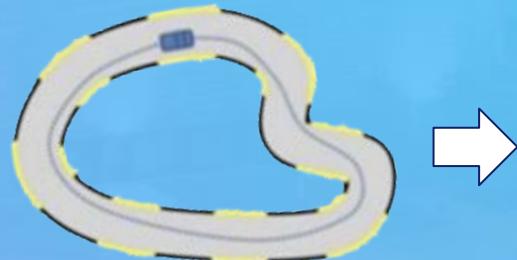
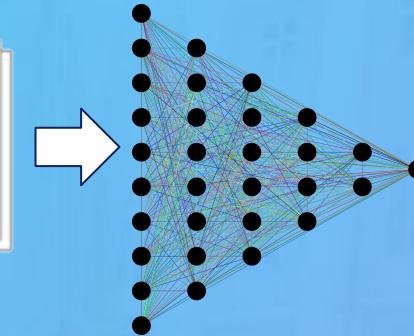
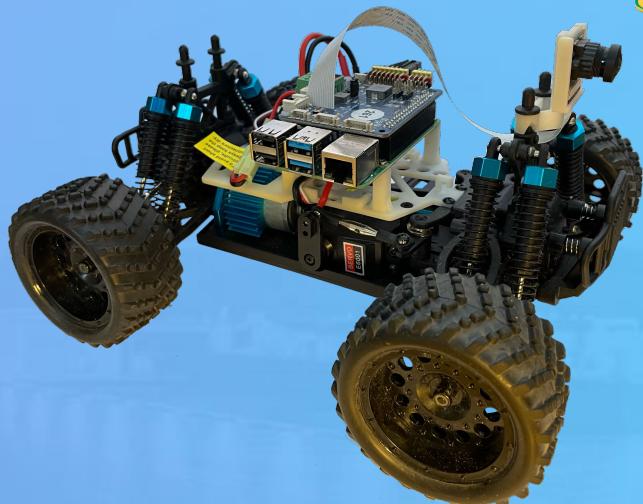


IMAGE +
JSON FILE
Steering Angle: 1.0
Throttle: 0.5
Milliseconds: 44522

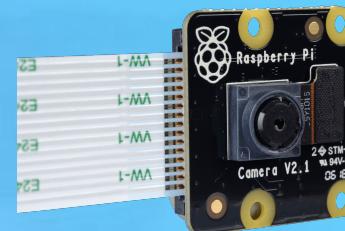


CNN Training

Method



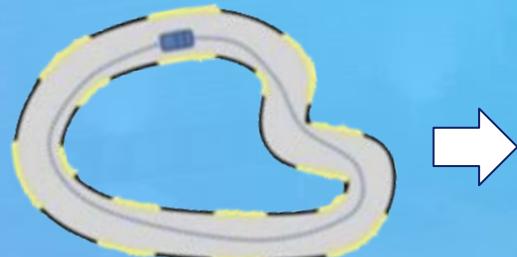
RaspberryPi



20 images
per second

Imitation Learning

TAL
TECH



Expert Demonstration

Training CNN models

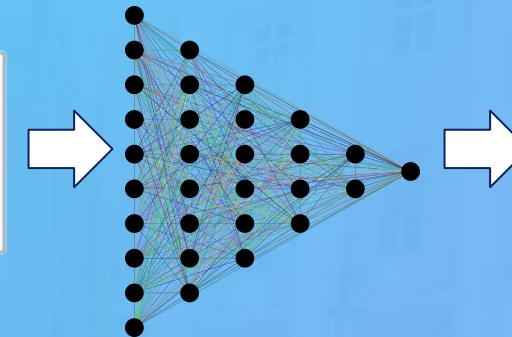
ONLY in lower-lights conditions



Dataset Cleaning



CNN Training



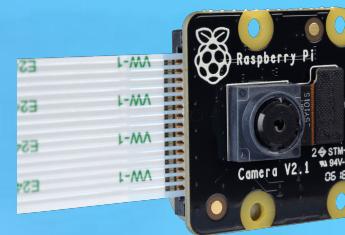
Model Output



Method



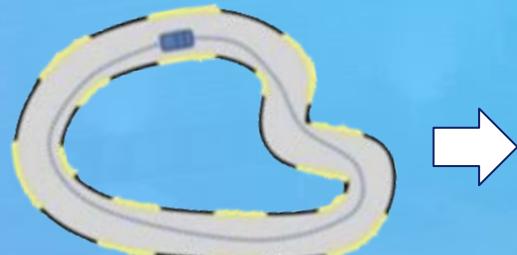
RaspberryPi



20 images per second

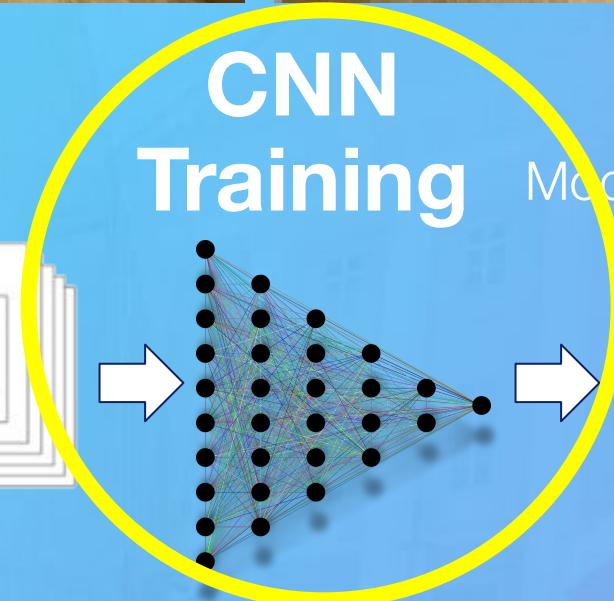
Expert Demonstration

Imitation Learning



Training CNN models

ONLY in lower-lights conditions



Dataset Cleaning



CNN
Training

Model Output



TAL
TECH

Method

Collect dataset
(only in lower-lights)

Training CNN models

Small

dataset
(~20 laps)



Method

Collect dataset
(only in lower-lights)

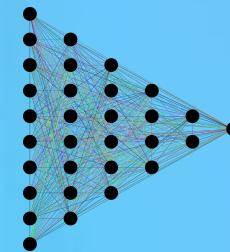
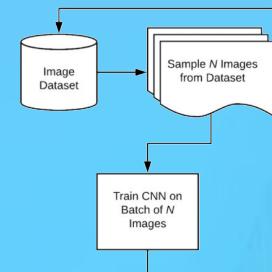
Small
dataset
(~20 laps)



Training CNN models

Method CNN training **CNN** output models

Standard



M-TS

Method

Collect dataset
(only in lower-lights)

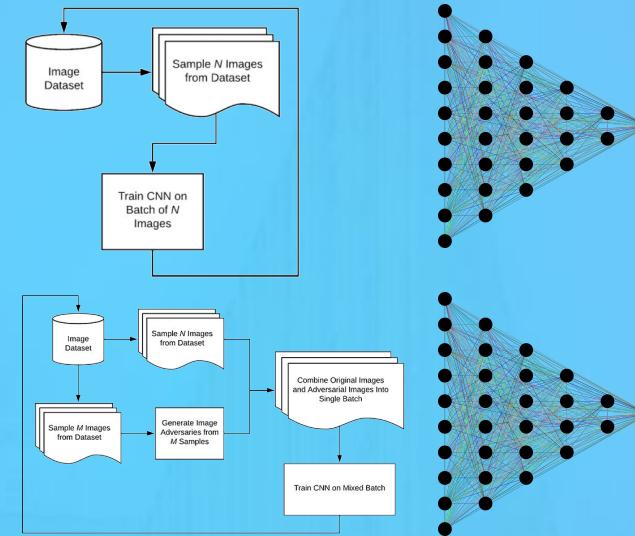
Small
dataset
(~20 laps)



Training CNN models

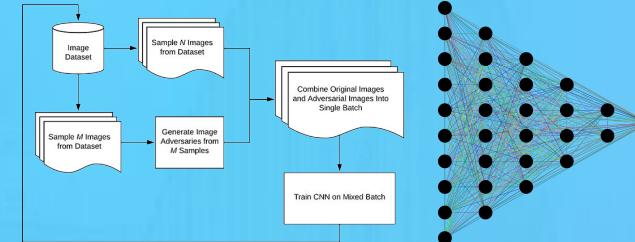
Method CNN training CNN output models

Standard



M-TS

Augmented



M-TSA

Method

Training CNN models

Collect dataset
(only in lower-lights)

Small
dataset
(~20 laps)



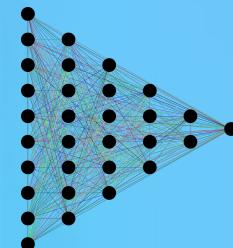
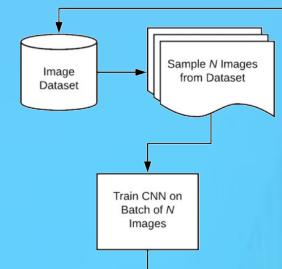
Large
dataset
(~40 laps)



**TAL
TECH**

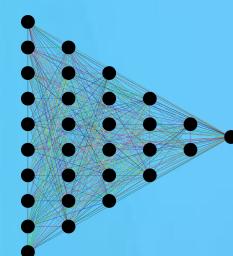
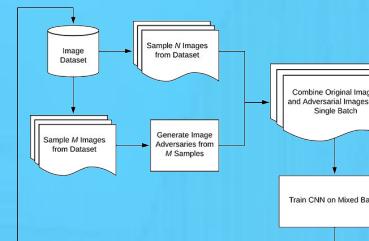
Method CNN training **CNN** output models

Standard



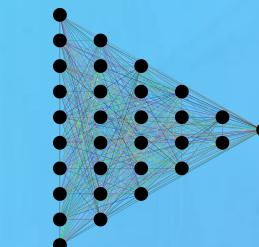
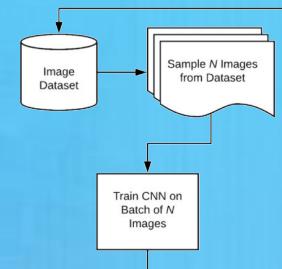
M-TS

Augmented



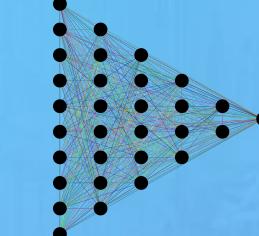
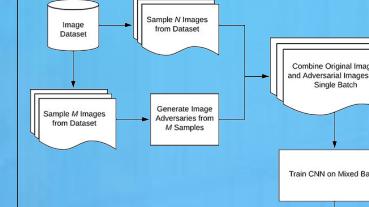
M-TSA

Standard



M-TL

Augmented



M-TLA

Method

Evaluating CNN models

Run 2 laps in the following conditions

Lower-lights

L



Higher-lights

H



M-TS



M-TL



M-TSA



M-TLA



Method

Evaluating CNN models

Run 2 laps in the following conditions

Lower-lights

L



Higher-lights

H



M-TS



M-TL



M-TSA



M-TLA



Method

Evaluating CNN models

Run 2 laps in the following conditions

Lower-lights

L



Higher-lights

H



M-TS



M-TL



M-TSA



M-TLA



Method

Evaluating CNN models

Run 2 laps in the following conditions

Lower-lights

L



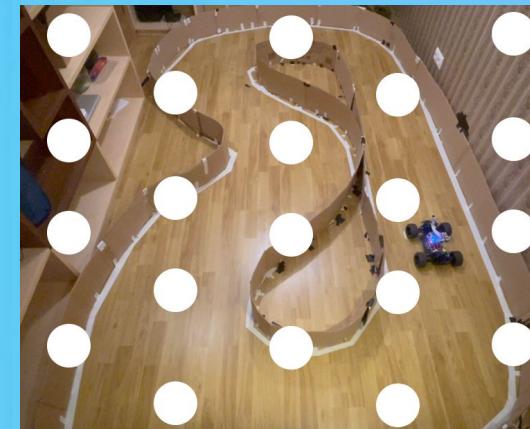
Higher-lights

H

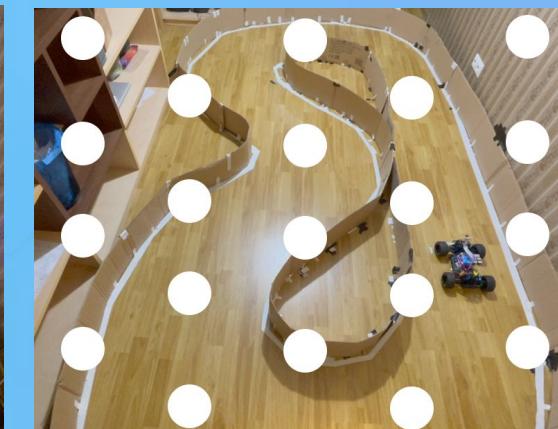


Lower-lights corrupted Higher-lights corrupted

LC



HC



M-TS



M-TL



M-TSA



M-TLA



Method

Evaluating CNN models

Run 2 laps in the following conditions

Lower-lights

L



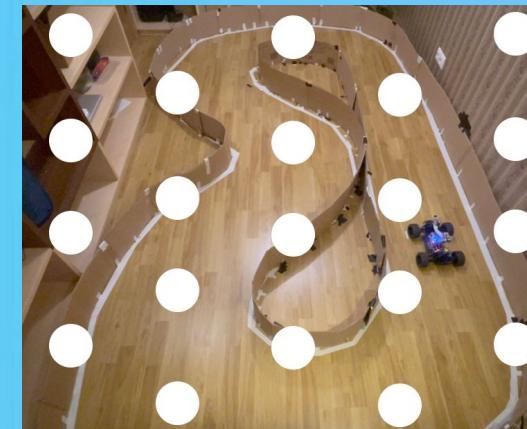
Higher-lights

H

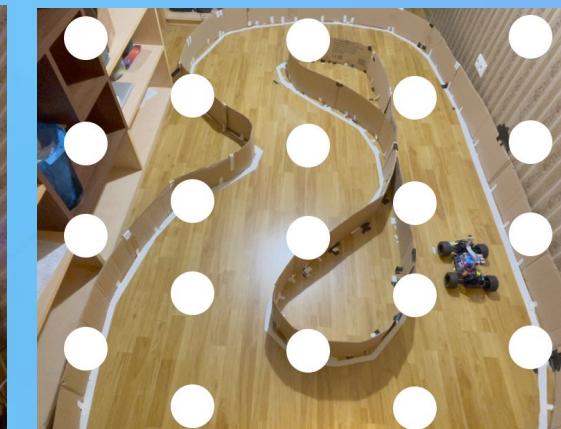


Lower-lights corrupted Higher-lights corrupted

LC



HC



- M-TS →
- M-TL →
- M-TSA →
- M-TLA →



?

?

Method

Evaluating CNN models

Run 2 laps in the following conditions

Lower-lights

L



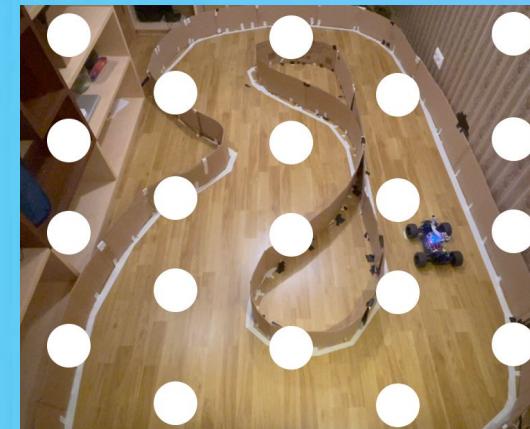
Higher-lights

H

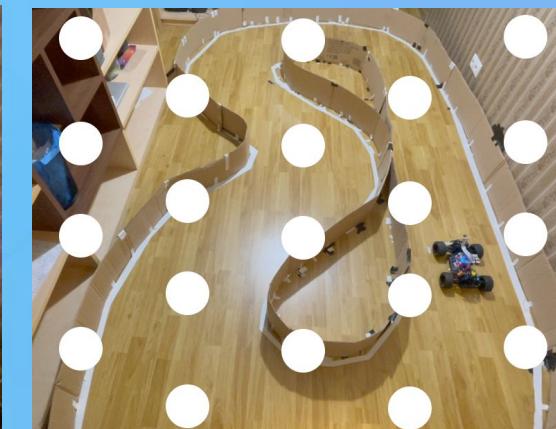


Lower-lights corrupted

LC



HC



M-TS



M-TL



M-TSA



M-TLA



$$4 \text{ evaluated in } 4 = 16$$

Models Conditions Evaluations

Method

Evaluating CNN models

Lighting conditions

Lower-lights

L



Higher-lights

H



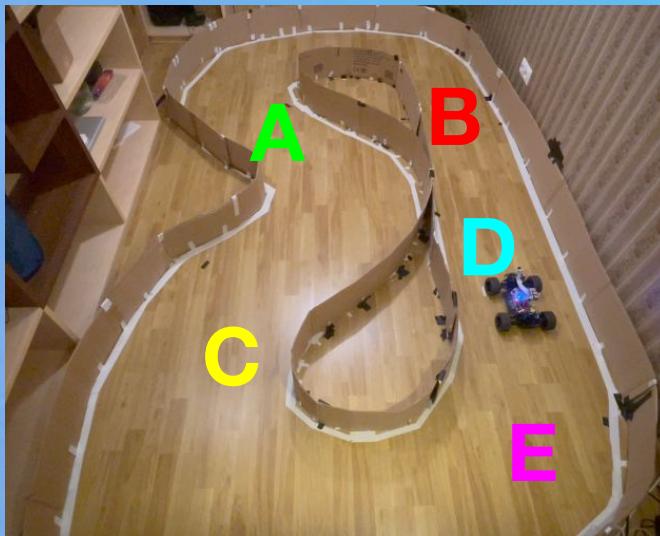
Method

Evaluating CNN models

Lighting conditions

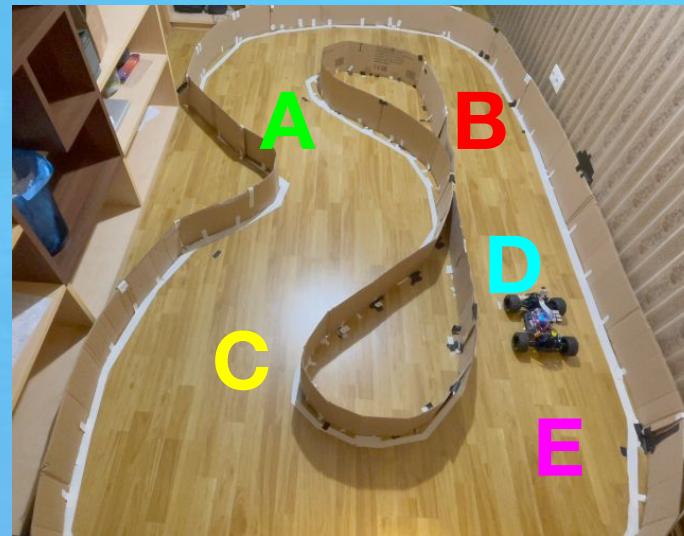
Lower-lights

L



Higher-lights

H



Location	Lower-lights L	Higher-lights H
A	21	110
B	19	86
C	19	75
D	19	75
E	20	60
Average	19.6 lux	81.2 lux

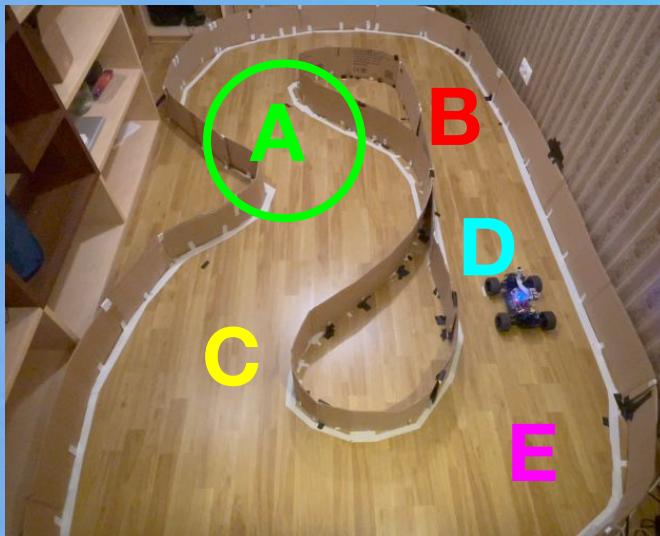
Method

Evaluating CNN models

Lighting conditions

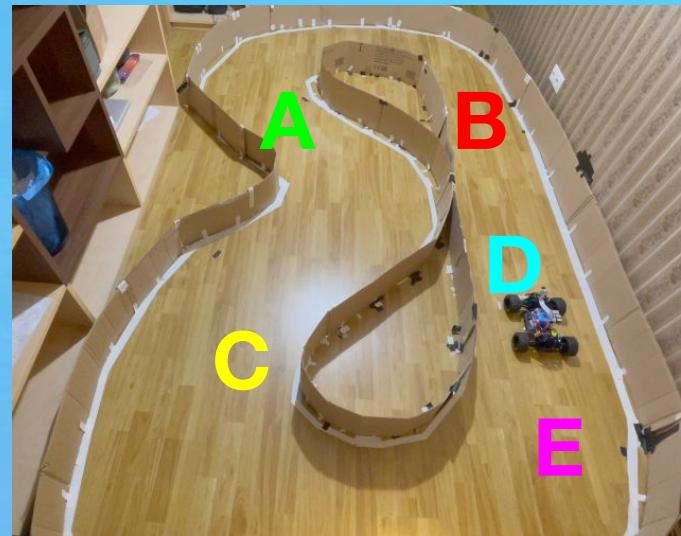
Lower-lights

L



Higher-lights

H



Location	Lower-lights	Higher-lights
A	21	110
B	19	86
C	19	75
D	19	75
E	20	60
Average	19.6 lux	81.2 lux

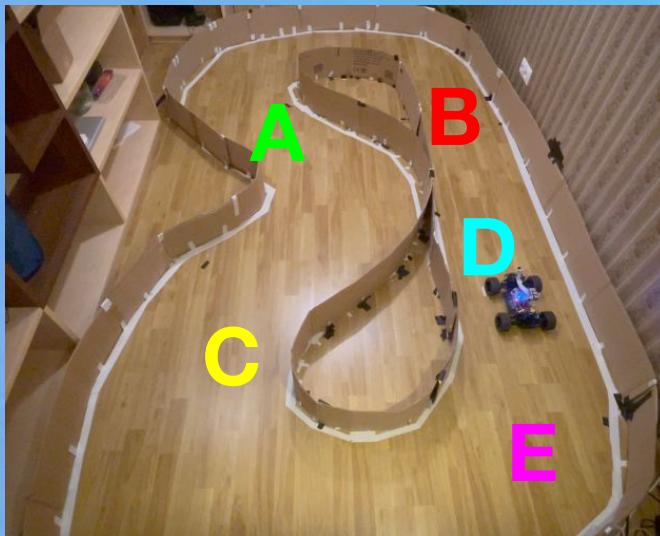
Method

Evaluating CNN models

Lighting conditions

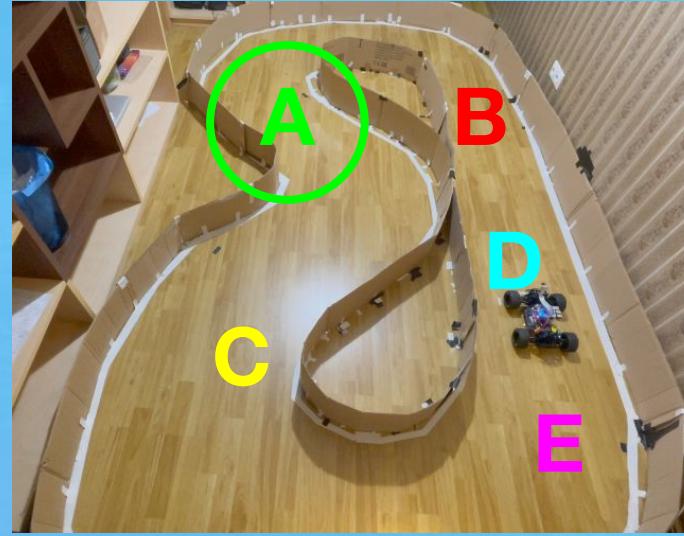
Lower-lights

L



Higher-lights

H



Location	Lower-lights	Higher-lights
A	21	110
B	19	86
C	19	75
D	19	75
E	20	60
Average	19.6 lux	81.2 lux

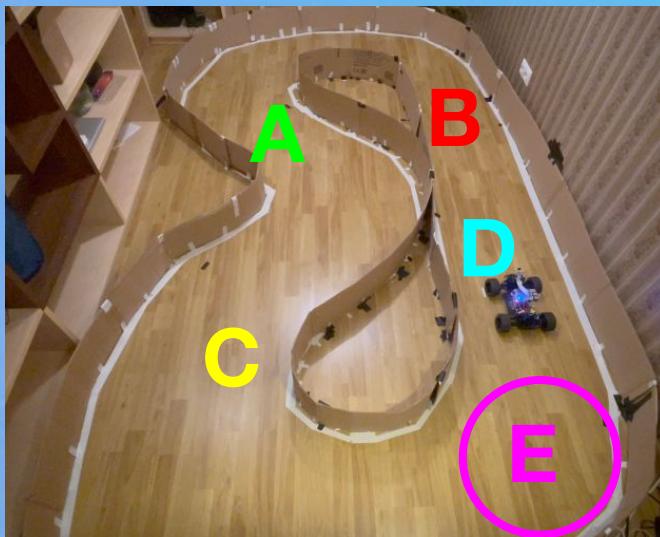
Method

Evaluating CNN models

Lighting conditions

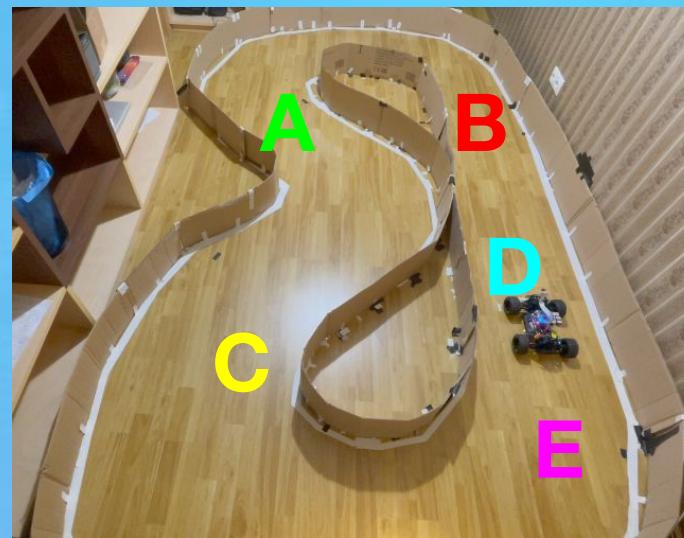
Lower-lights

L



Higher-lights

H



Location	Lower-lights L	Higher-lights H
A	21	110
B	19	86
C	19	75
D	19	75
E	20	60
Average	19.6 lux	81.2 lux

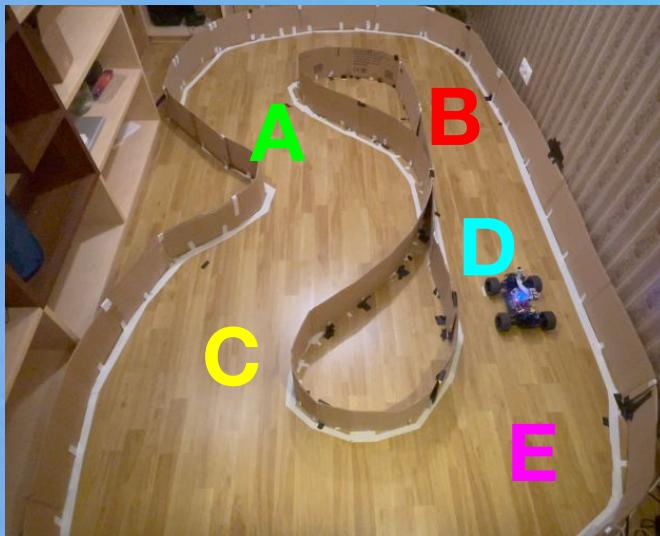
Method

Evaluating CNN models

Lighting conditions

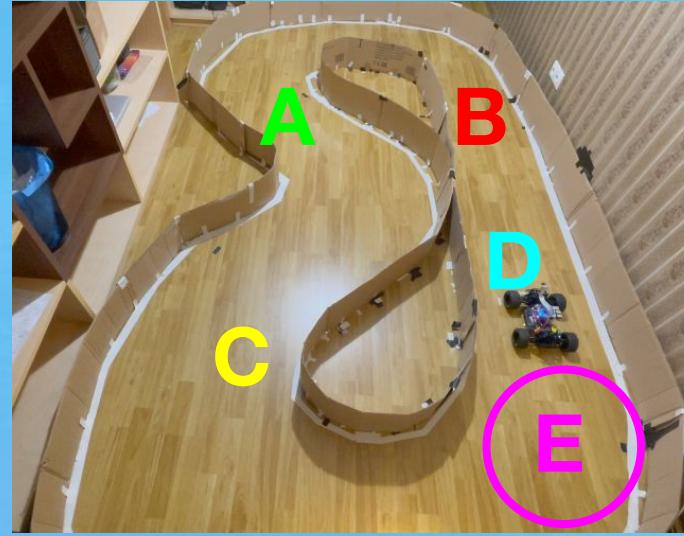
Lower-lights

L



Higher-lights

H



Location	Lower-lights	Higher-lights
A	21	110
B	19	86
C	19	75
D	19	75
E	20	60
Average	19.6 lux	81.2 lux

Evaluating conditions (continued)

Independent variable	Values
Laps	2
Lights	low, high
Vision	uncorrupted, corrupted
Dependent variables	Values
Collisions	[0...n]

Results

Collisions in 2 laps

👁 Seen conditions

Lower-lights

L



Higher-lights

H



❗ Never-seen-before conditions

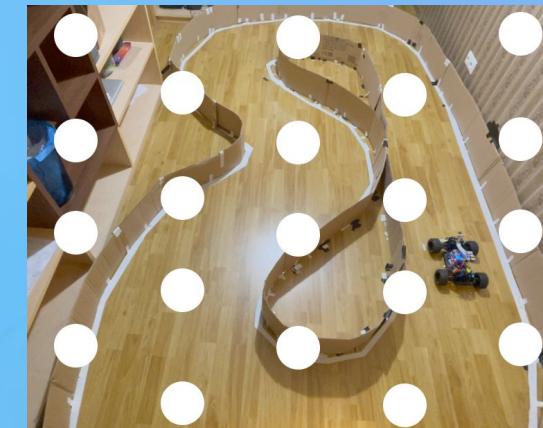
Lower-lights corrupted

LC



Higher-lights corrupted

HC



M-TS

M-TL

M-TSA

M-TLA

TAL
TECH

Results

Collisions in 2 laps

👁 Seen conditions

Lower-lights

L



Higher-lights

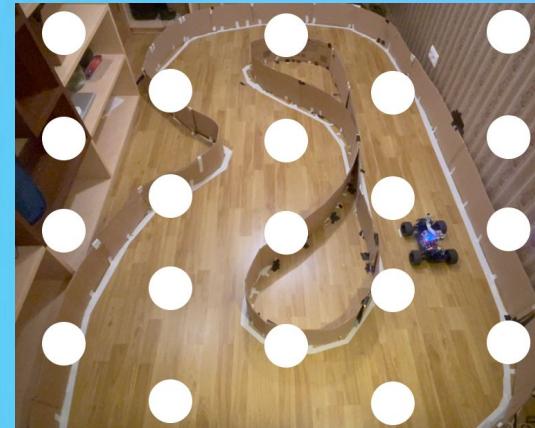
H



❗ Never-seen-before conditions

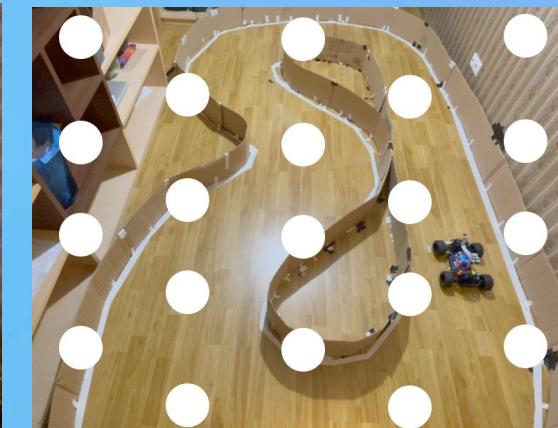
Lower-lights corrupted

LC



Higher-lights corrupted

HC



M-TS



0

M-TL



0

M-TSA

M-TLA

Results

Collisions in 2 laps

👁 Seen conditions

Lower-lights

L



Higher-lights

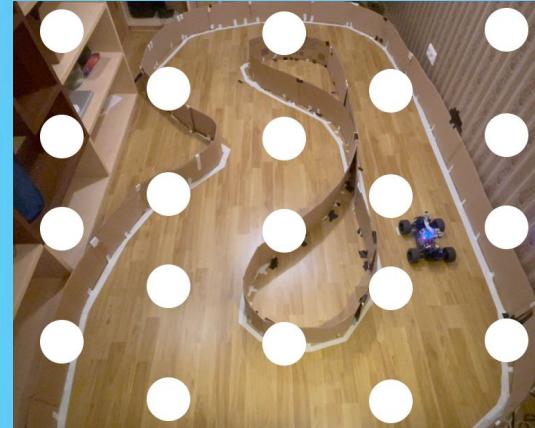
H



❗ Never-seen-before conditions

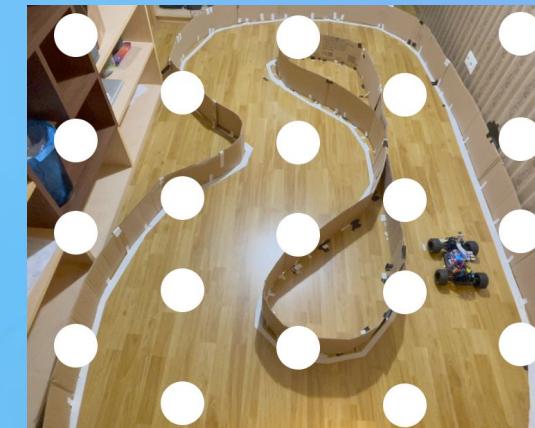
Lower-lights corrupted

LC



Higher-lights corrupted

HC



M-TS



0

M-TL



0

M-TSA



0

M-TLA



0

Results

Collisions in 2 laps

👁 Seen conditions

Lower-lights

L



Higher-lights

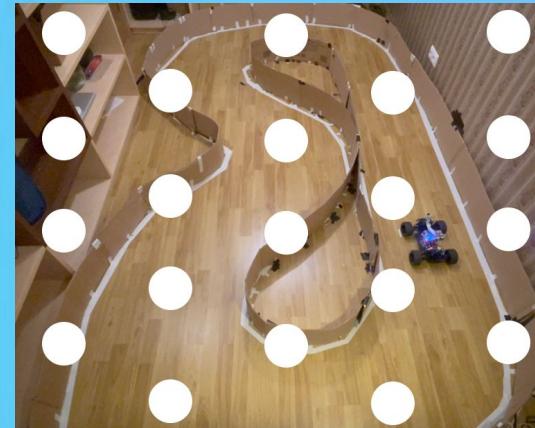
H



❗ Never-seen-before conditions

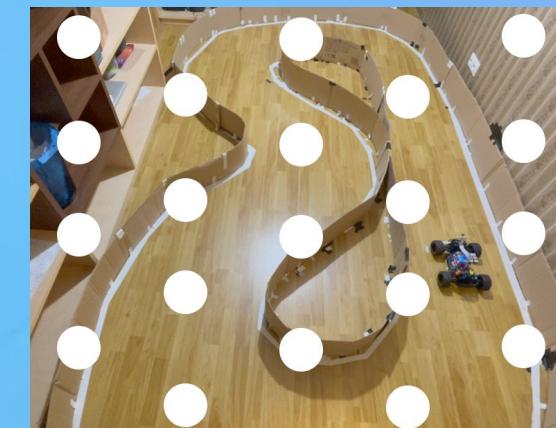
Lower-lights corrupted

LC



Higher-lights corrupted

HC



M-TS



0



5

M-TL



0



4

M-TSA



0

M-TLA



0

Results

Collisions in 2 laps

👁 Seen conditions

Lower-lights

L



Higher-lights

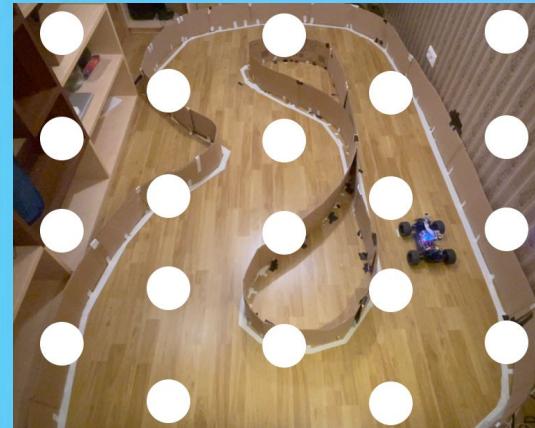
H



❗ Never-seen-before conditions

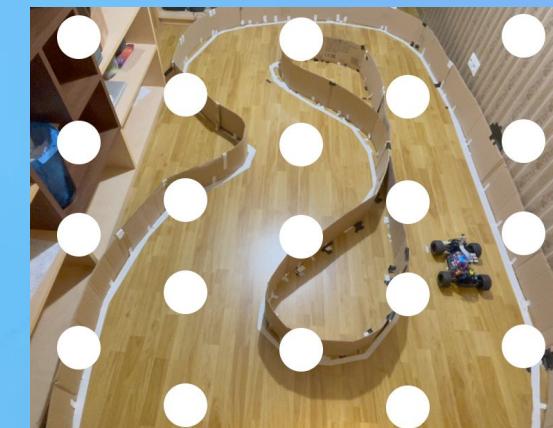
Lower-lights corrupted

LC



Higher-lights corrupted

HC



M-TS

✓ 0

M-TL

✓ 0

M-TSA

✓ 0

M-TLA

✓ 0

✗ 5
✗ 4



Results

Collisions in 2 laps

👁 Seen conditions

Lower-lights

L



Higher-lights

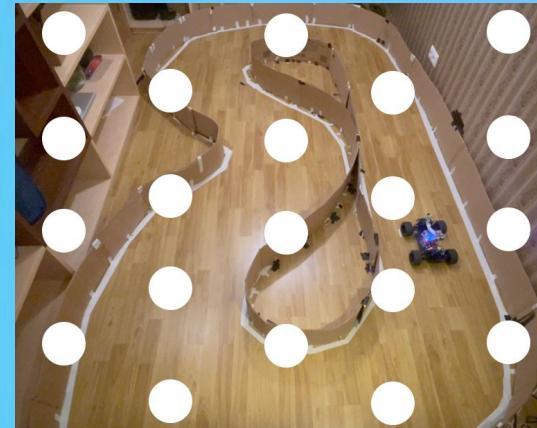
H



❗ Never-seen-before conditions

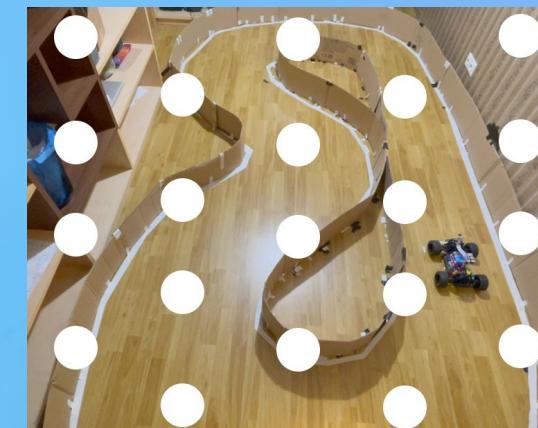
Lower-lights corrupted

LC



Higher-lights corrupted

HC



M-TS



0



5

M-TL



0



4

M-TSA



0



0



M-TLA



0



0

Results

Collisions in 2 laps

👁 Seen conditions

Lower-lights

L



Higher-lights

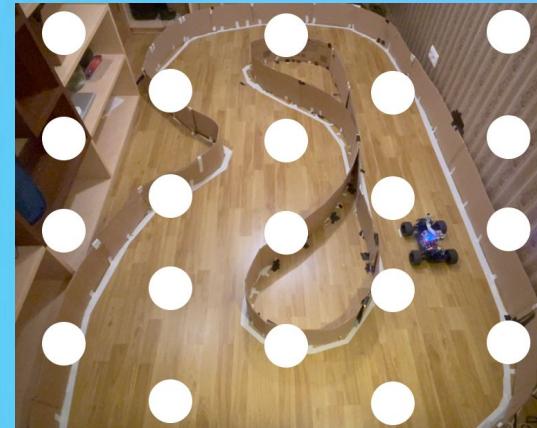
H



❗ Never-seen-before conditions

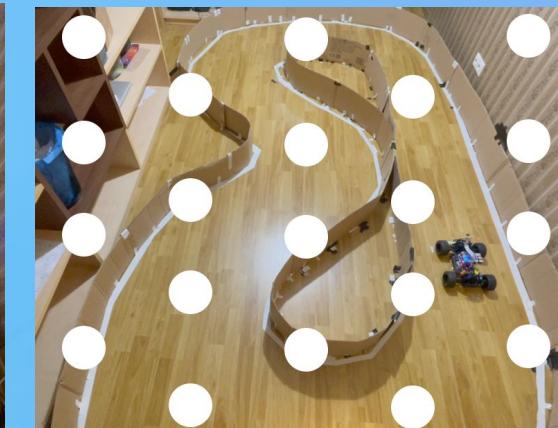
Lower-lights corrupted

LC



Higher-lights corrupted

HC



M-TS



0



5

M-TL



0



4

M-TSA



0



0



2



5

M-TLA



0



0



3



5

Results

Collisions in 2 laps

👁 Seen conditions

Lower-lights

L



Higher-lights

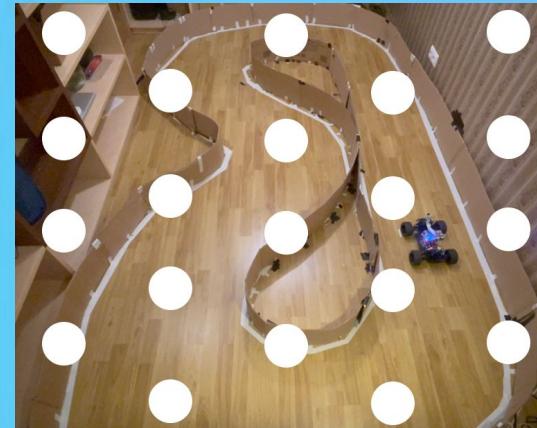
H



❗ Never-seen-before conditions

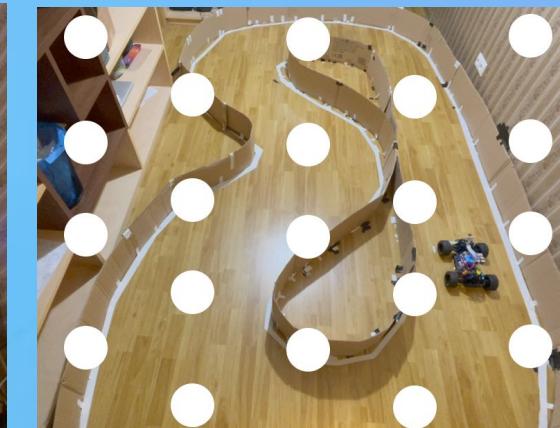
Lower-lights corrupted

LC



Higher-lights corrupted

HC



M-TS

✓ 0

✗ 5

✗ 12

✗ 7

M-TL

✓ 0

✗ 4

✗ 10

✗ 9

M-TSA

✓ 0

✓ 0

✗ 2

✗ 5

M-TLA

✓ 0

✓ 0

✗ 3

✗ 5

Results

Collisions in 2 laps

👁 Seen conditions

Lower-lights

L



Higher-lights

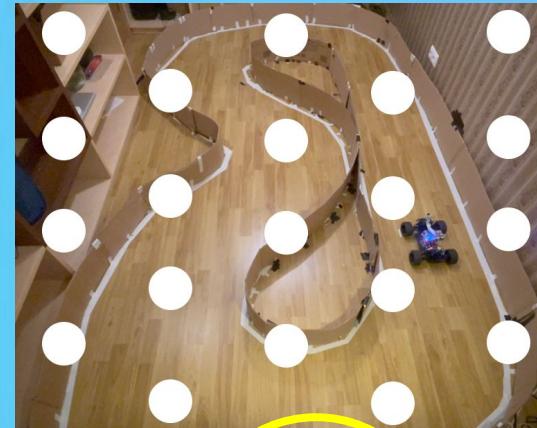
H



🤔 Never-seen-before conditions

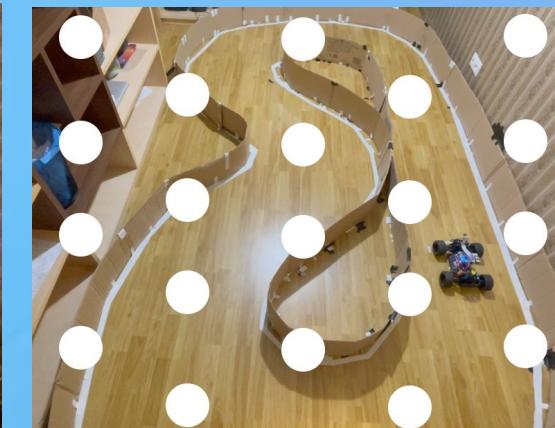
Lower-lights corrupted

LC



Higher-lights corrupted

HC



M-TS

✓ 0

✗ 5

✗ 12

✗ 7

M-TL

✓ 0

✗ 4

✗ 10

✗ 9

M-TSA

✓ 0

✓ 0

✗ 2

✗ 5

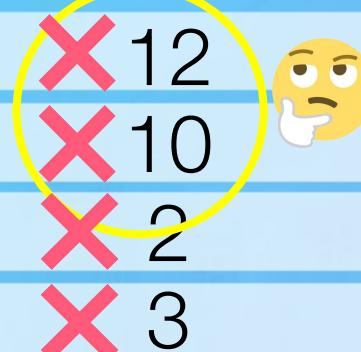
M-TLA

✓ 0

✓ 0

✗ 3

✗ 5



Answer to the Research Question

“Can adversarial defense training methods improve the neural network generalization skills to unseen lighting conditions?”

Answer to the Research Question

“Can adversarial defense training methods improve the neural network generalization skills to unseen lighting conditions?”

Yes, it can.

Answer to the Research Question

“Can adversarial defense training methods improve the neural network generalization skills to unseen lighting conditions?”

Yes, it can.

Standard method

- ✓ Seen lower-lights
- ✗ Never seen before higher-lights

Augmented training method

- ✓ Seen lower-lights
- ✓ Never seen before higher-lights

Answer to the Research Question

“Can adversarial defense training methods improve the neural network generalization skills to unseen lighting conditions?”

Yes, it can.

Standard method

- ✓ Seen lower-lights
- ✗ Never seen before higher-lights
- ✗ Vulnerable to adversarial attack

Augmented training method

- ✓ Seen lower-lights
- ✓ Never seen before higher-lights
-  Reduced collisions under attack

Future Work



Measure differences with training dataset and evaluation dataset



Image augmentations



More diverse selection of attack techniques



Use datasets constructed only from corrupted images



Larger volume of samples

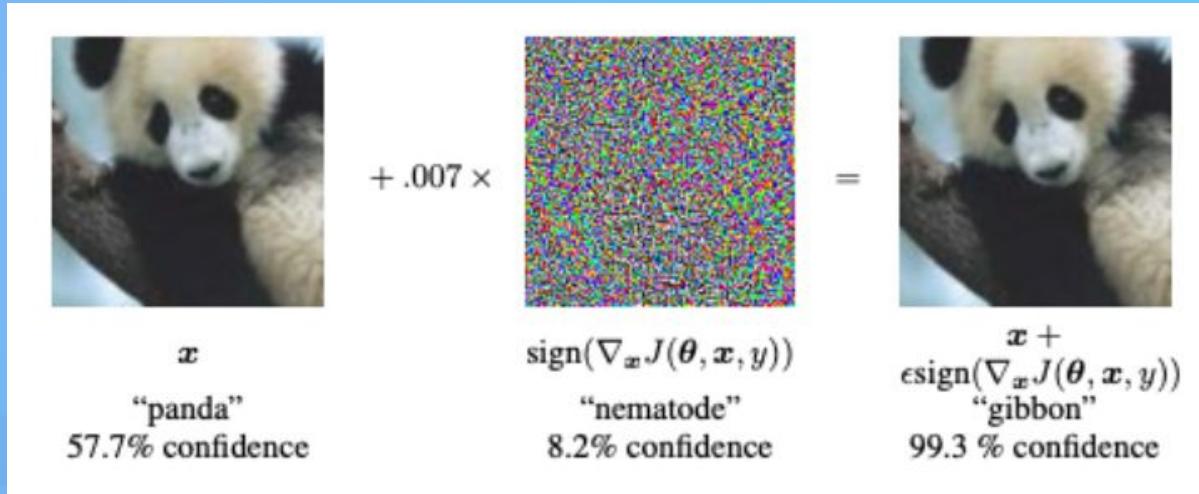
Using Adversarial Defense Methods to Improve the Performance
of Deep-Neural-Network-Controlled Automatic Driving Systems

Mike Camara

Questions?

Backup Slides

Machine Learning Adversarial Attacks



Explaining and Harnessing Adversarial Examples (Goodfellow, Shlens, & Szegedy, 2015)

White Box

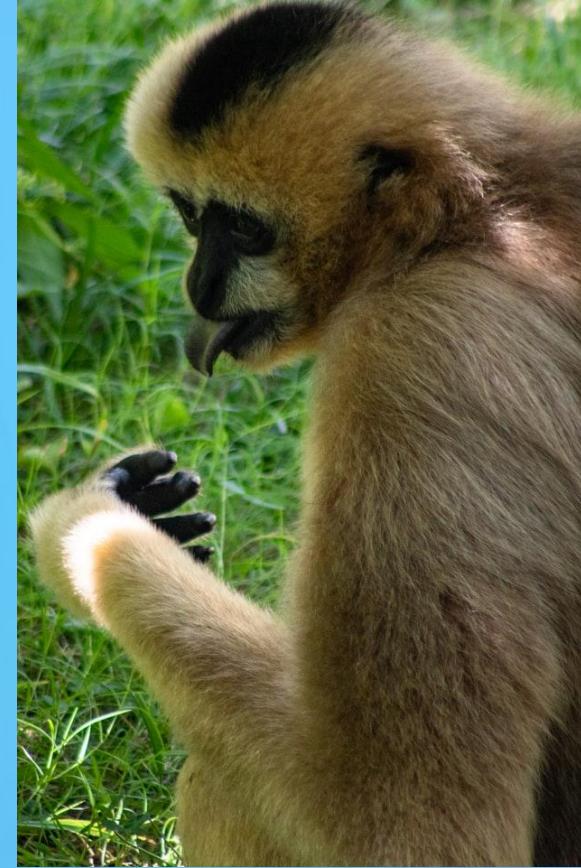
Evasion modify input to affect model (FGSM)

Black Box

Poisoning modify training data to add vulnerabilities

Inference learn confidential private data

Extraction theft of a proprietary model



📸 A Gibbon by @jcotten

Adversarial Robustness Toolbox - ART (IBM Think, 2021)

Attacks

Evasion

Fast Gradient Sign Method (FGSM) – white-box

DeepFool – white-box

Virtual Adversarial Method – white-box

Carlini & Wagner L2 and Linf – white-box

Basic Iterative Method – white-box

Jacobian Saliency Map – white-box

Universal Perturbation – white-box

Decision Tree Attack – white-box

Projected Gradient Descent (PGD) – white-box

NewtonFool – white-box

Elastic Net – white-box

Spatial Transformation – black-box

Query-efficient Black-box – black-box

Zeroth-Order Optimization (ZOO) – black-box

Adversarial Patch – white-box

Boundary Attack – black-box

High Confidence Low Uncertainty (HCLU) – white-box

HopSkipJump – black-box

AutoAttack

AutoPGD – white-box

Brendel&Bethge – white-box

Feature Adversaries – white-box

Shadow Attack – white-box

Simple Black-box Attack (SimBA) – black-box

Wasserstein Attack – white-box

Geometric Decision-based Attack (GeoDA) – black-box

Poisoning

Poisoning Attack on SVM

Adversarial Embedding

Bullseye Polytope

Clean-label Backdoor

Feature Collision

Extraction

Functionally Equivalent Extraction

Copycat CNN

Knockoff Nets

Inference

Membership inference

Attribute inference

Model inversion

Database reconstruction

Defences:

Adversarial Training

Adversarial Training

Adversarial Training Madry-PGD

Fast is Better than Free

Preprocessing

Thermometer encoding

Total variance minimization

PixelDefend

Gaussian data augmentation

Feature squeezing

Spatial smoothing

JPEG compression

Label smoothing

Postprocessing

Reverse Sigmoid

Random Noise

Class Labels

High Confidence

Rounding

Certification

Randomized Smoothing

Metrics

CLEVER

Loss sensitivity

Empirical robustness

Verification

Clique Method Robustness Verification

Attack and Fault Injection in Self-Driving Agents on the CARLA Simulator

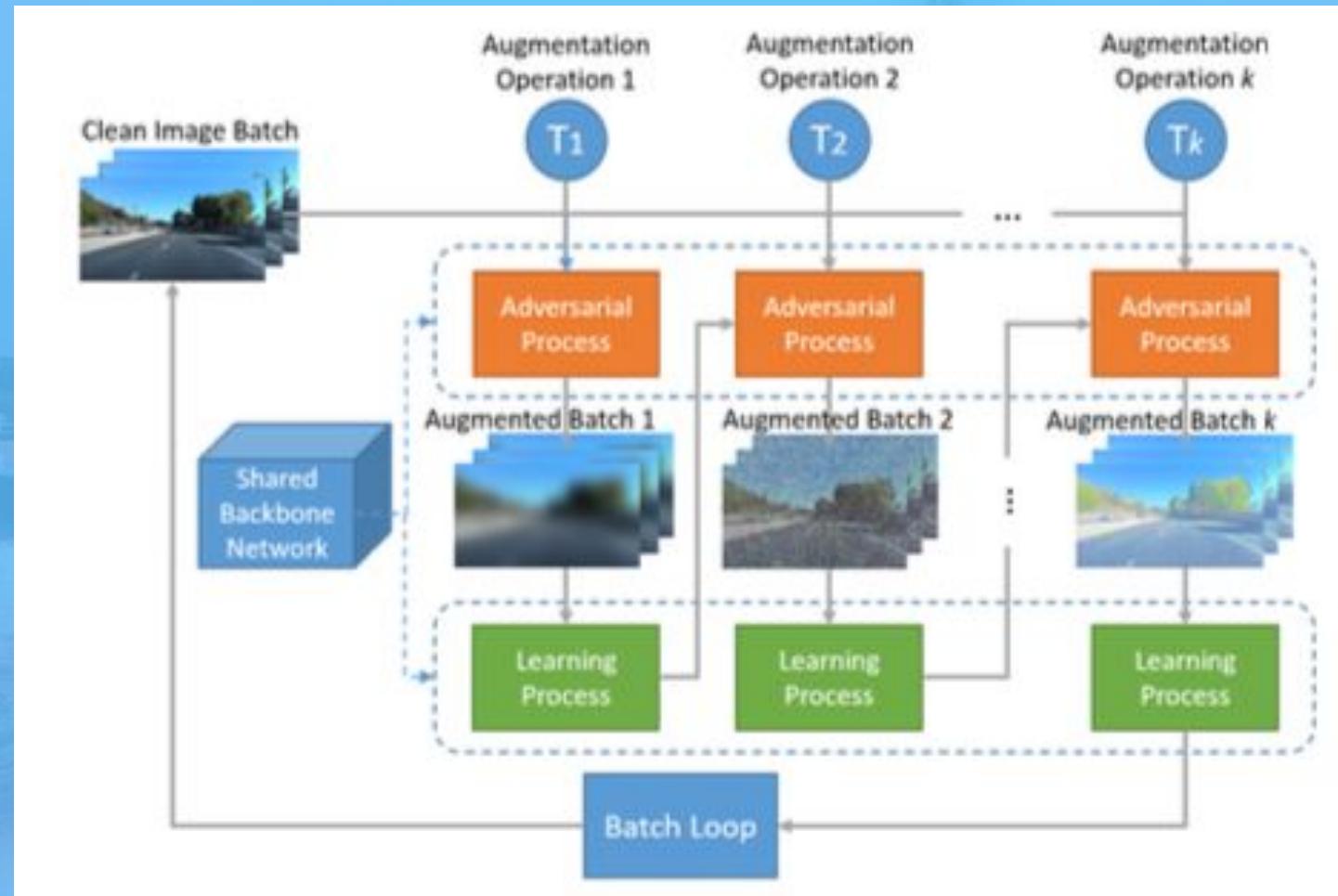
(Piazzesi, Hong, & Ceccarelli, 2021)



Measured impact of a range of adversarial attacks on the performance of Automatic Driving CNN's deployed in CARLA Simulator environment react to several different adversarial attacks by analyzing the rate of collisions in various circuits in a specific set of predefined conditions.

Adversarial Differentiable Data Augmentation for Autonomous Systems

(M. Shu, Y. Shen, M. C. Lin, and T. Goldstein, 2021)



"Adversarially training against image corruptions it results having CNN's with higher levels of robustness and accuracy for a range of settings as compared to baseline and state-of-the-art augmentation methods."

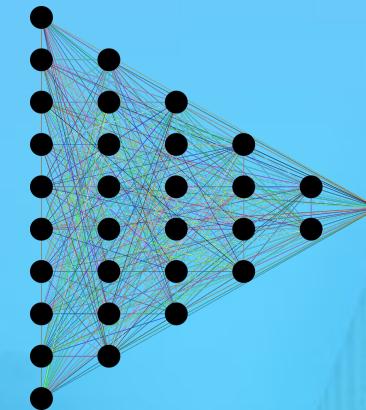
Mixing Normal Images and Adversarial Images when Training CNNs

(A. Rosebrock, 2021)

```
normal testing images:  
loss: 0.0477, acc: 0.9891
```

```
generating adversarial examples with FGSM...
```

```
adversarial testing images:  
loss: 14.0658, acc: 0.0188
```

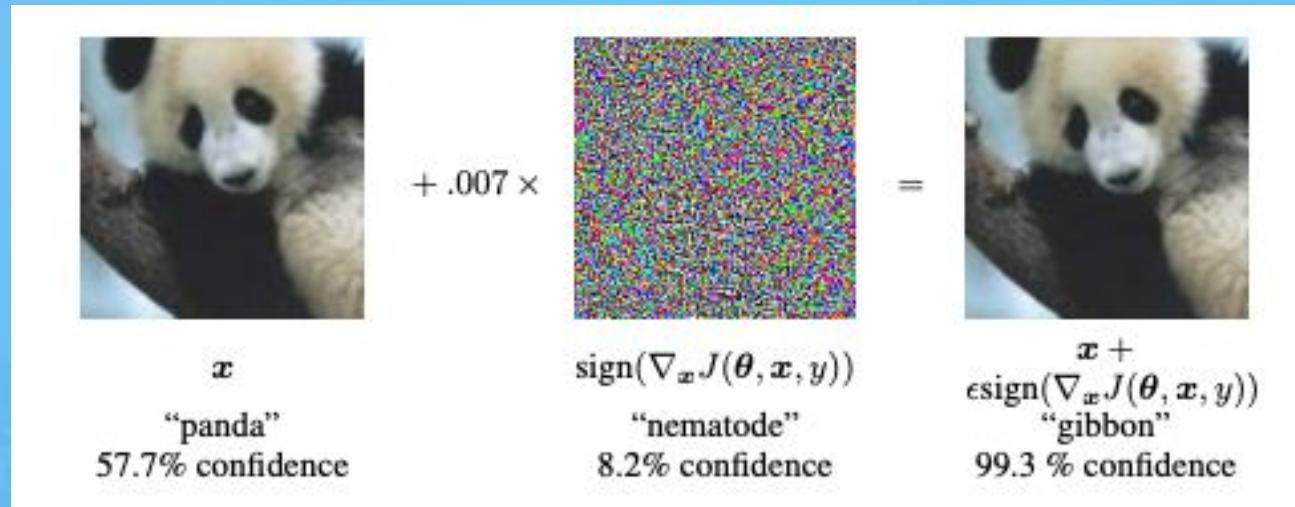


```
normal testing images *after* fine-tuning:  
loss: 0.0315, acc: 0.9906
```

```
adversarial images *after* fine-tuning:  
loss: 0.1190, acc: 0.9641
```

"Mixed batch training will make your model become more robust and generalize better" - Rosebrock

Fast Gradient Sign Method (FGSM)



(Goodfellow, Shlens, & Szegedy, 2015)



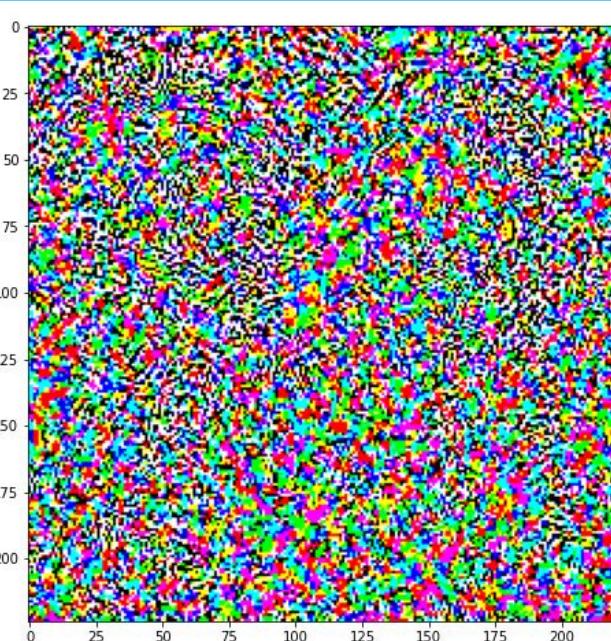
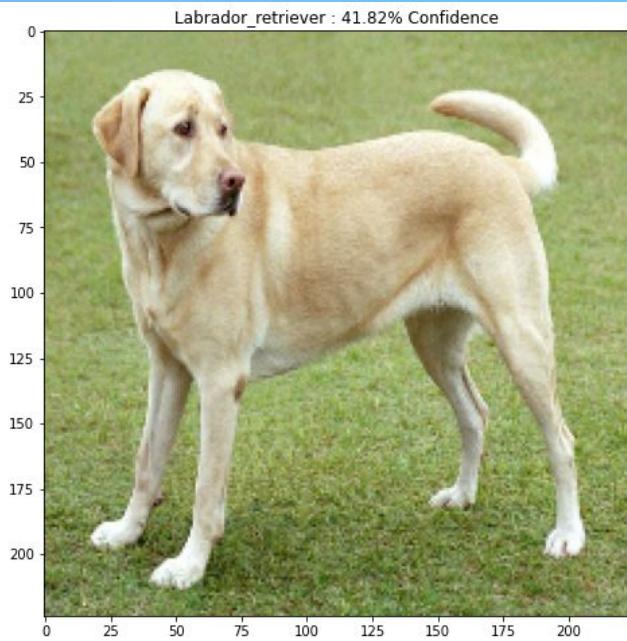
$$\epsilon \\ 0.8$$

Normal

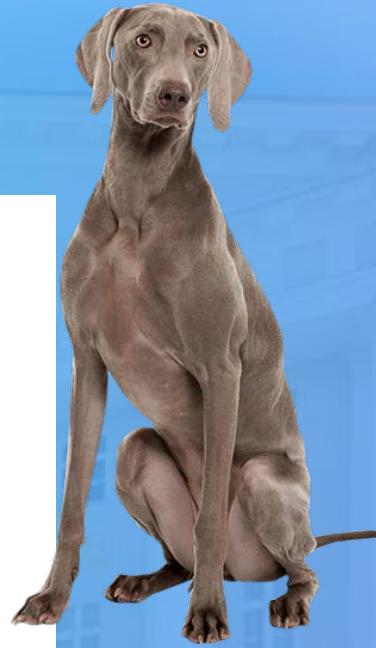
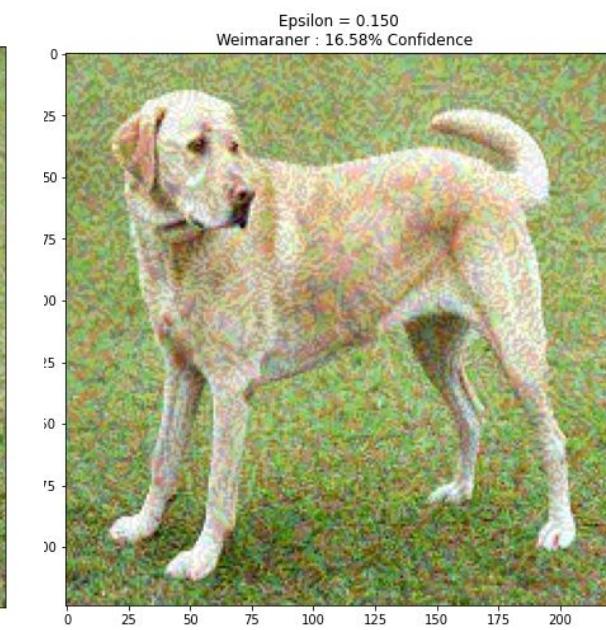
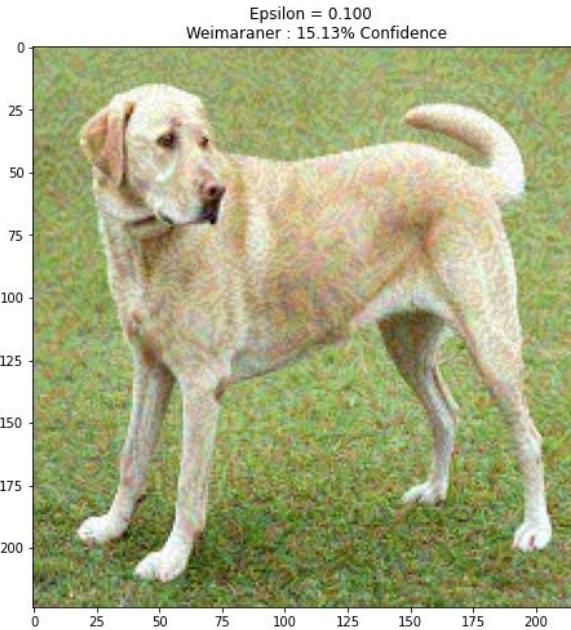


Corrupted

Adversarial Example Using FGSM (TensorFlow, 2022)



MobileNetV2 model and
the ImageNet class names

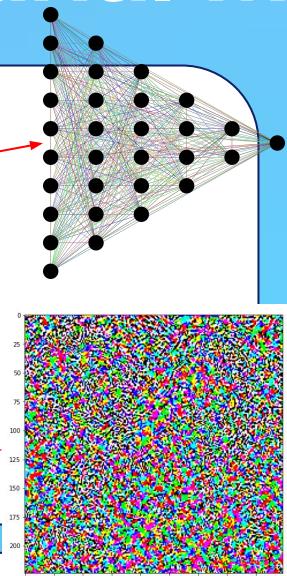


Adversarial Machine Learning

```
def adversarial_pattern(self, model, image, label):
    image = tensorflow.cast(image, tensorflow.float32)
    with tensorflow.GradientTape() as tape:
        tape.watch(image)
        prediction = model(image)
        loss = tensorflow.keras.losses.MSE(label, prediction)

    gradient = tape.gradient(loss, image)
    signed_grad = tensorflow.sign(gradient)

    return signed_grad
```



```
new_data_folder = Tub(os.path.join(base_path, 'corrupted_images/'),
                      inputs=['cam/image_array', 'user/angle', 'user/throttle'],
                      types=['image_array', 'float', 'float'])

for key, record in enumerate(tub1):
    t_record = TubRecord(config=cfg,
                          base_path=tub1.base_path,
                          underlying=record)
    img_arr = t_record.image(cached=False)
    record['user/angle'] = record['user/angle']
    record['user/throttle'] = record['user/throttle']

    image = img_arr.reshape((1,) + img_arr.shape)
    label_to_pass = model.predict(image)
    perturbation = self.adversarial_pattern(model, image, label_to_pass).numpy()
    perturb = ((perturbation[0]*0.5 + 0.5)*255)-50
    adv_img = np.clip(img_arr + (perturb*0.8), 0, 255)
    adv_img = adv_img.astype(int)
    record['cam/image_array'] = adv_img
    new_data_folder.write_record(record)
```

FGSM

$$\mathbf{x} + .007 \times \text{sign}(\nabla_{\mathbf{x}} J(\theta, \mathbf{x}, y)) = \mathbf{x} + \epsilon \text{sign}(\nabla_{\mathbf{x}} J(\theta, \mathbf{x}, y))$$

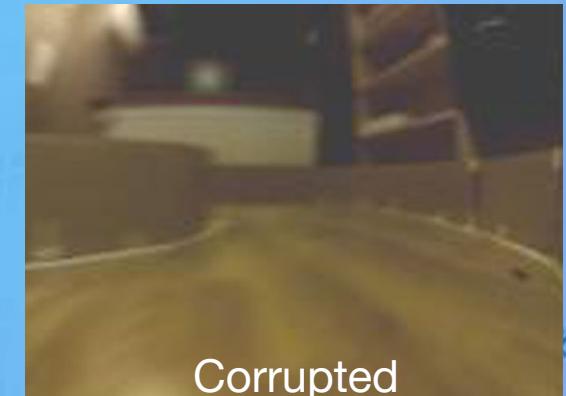
\mathbf{x} "panda" 57.7% confidence
+ .007 × sign($\nabla_{\mathbf{x}} J(\theta, \mathbf{x}, y)$) "nematode" 8.2% confidence
= $\mathbf{x} + \epsilon \text{sign}(\nabla_{\mathbf{x}} J(\theta, \mathbf{x}, y))$ "gibbon" 99.3 % confidence

(Goodfellow, Shlens, & Szegedy, 2015)

Batch of adversarial images generation



Normal

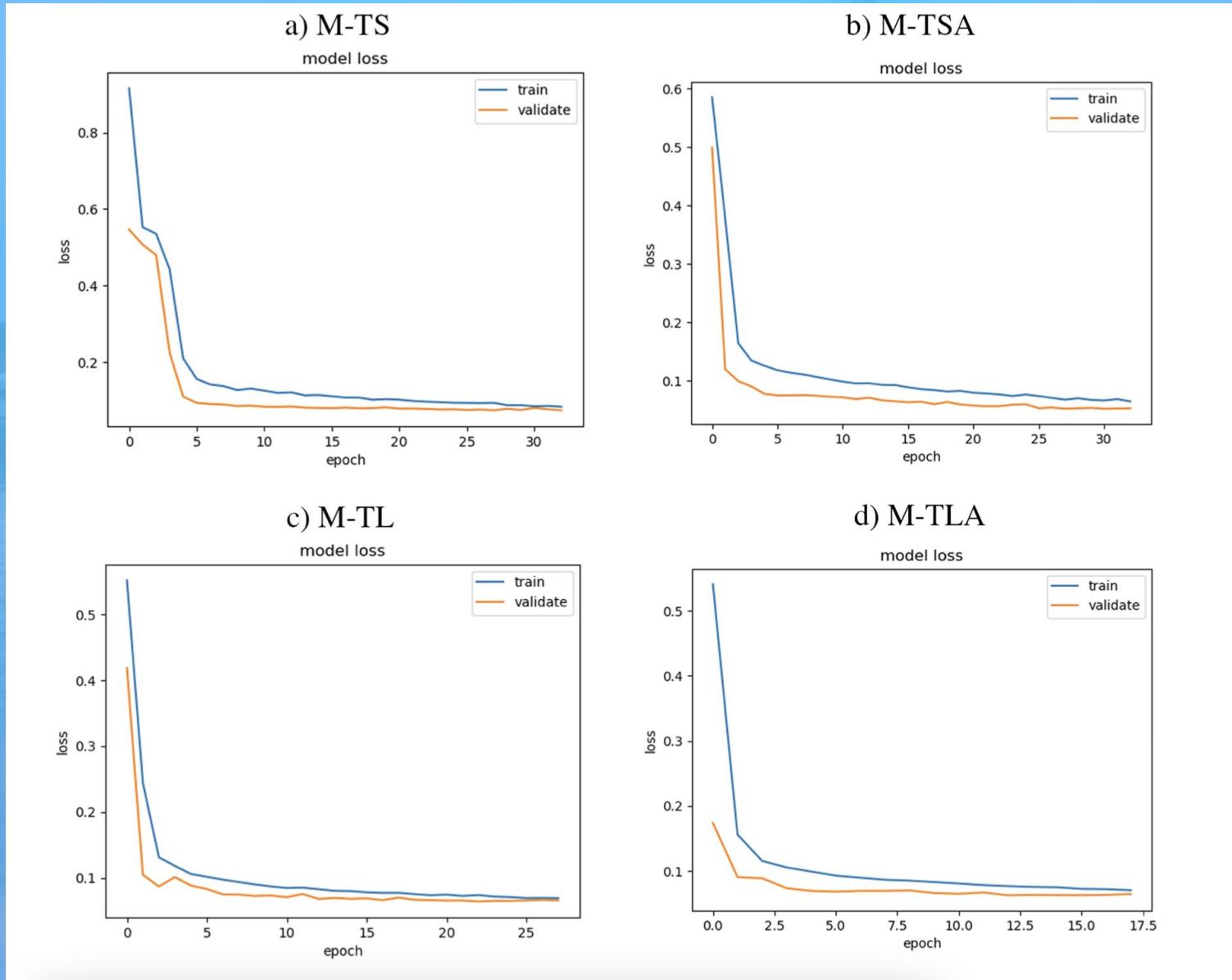


Corrupted

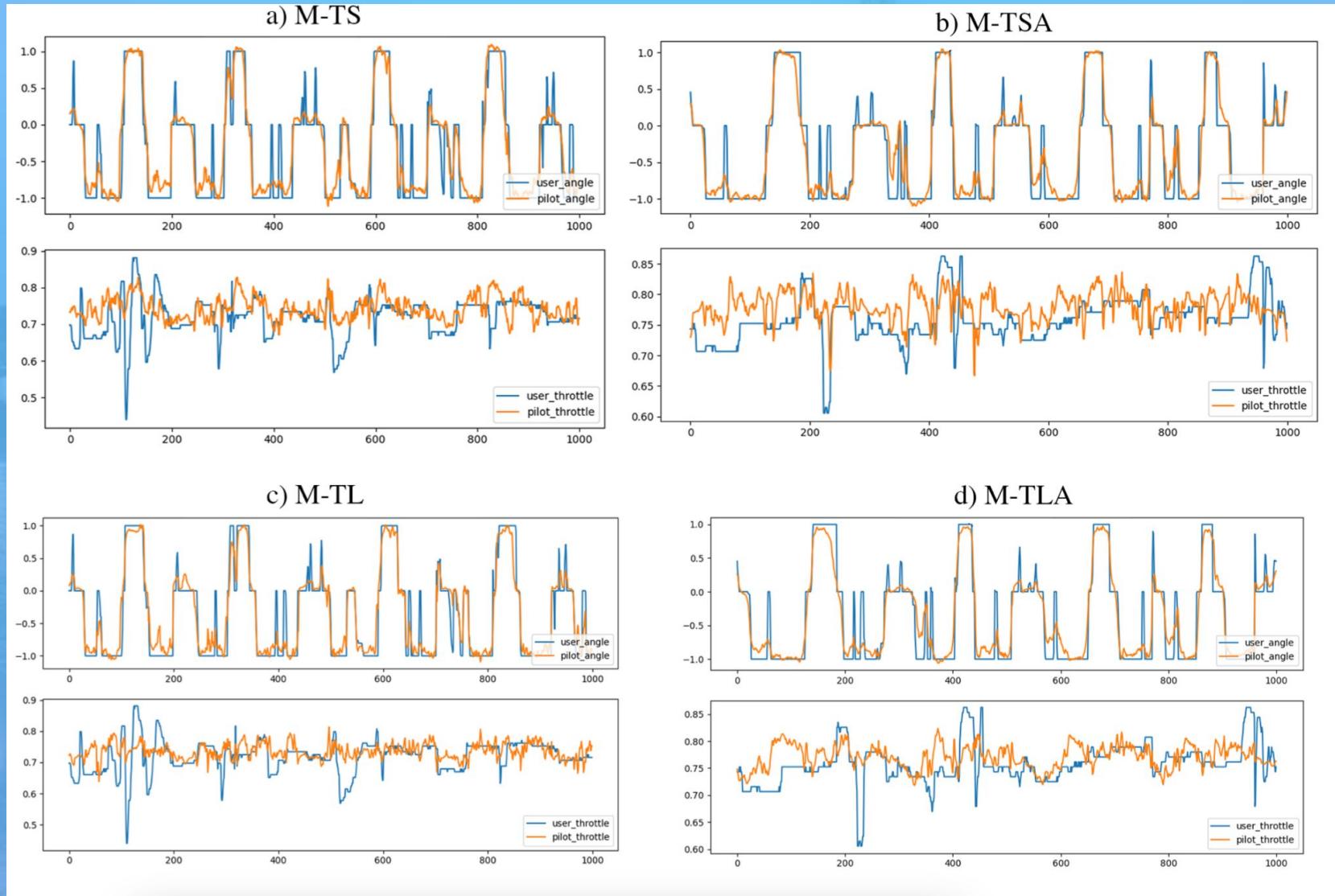
Data Cleaning

Dataset	Total images	Removed images
Small	5593	842
Large	10418	1491

Training



Performance



Evaluation

CNN Score - Quantifying the quality of predictions		
CNN Model	Prediction	RMSD
M-TS	Steering angle	0.24630
	Throttle	0.05860
M-TSA	Steering angle	0.22733
	Throttle	0.05193
M-TL	Steering angle	0.23190
	Throttle	0.04488
M-TLA	Steering angle	0.23527
	Throttle	0.04012

Society of Automotive Engineers (SAE) - Automation Levels



0

No
Automation

1

Driver
Assistance

2

Partial
Automation

3

Conditional
Automation

4

High
Automation

5

Full
Automation

End-to-end Pipeline



- Traffic-Aware Cruise Control
- Navigate on Autopilot
- Autosteer
- Traffic Light and Stop Sign Control
- Autopark
- Smart Summon

End-to-end Pipeline



✓ End-to-end driving a vehicle in traffic on public UK roads

End-to-end Pipeline



Carnegie
Mellon
University



ALVINN

The first neural network
powered automatic driving
system.

(Pomerleau, 1988)

End-to-end Pipeline



**Deep Learning Approach
to Automatic Driving Systems (ADS)**

Thank you





TAL
TECH

