

# Data Augmentation via Neural-Style-Transfer for Driver Distraction Recognition

Che-Tsung Lin<sup>1\*</sup>, Thomas Streubel<sup>2,3</sup>, Marco Dozza<sup>2</sup>, Giulio Bianchi Piccinini<sup>2</sup>, Christopher Zach<sup>1</sup>

<sup>1</sup> Department of Electrical Engineering, Chalmers University of Technology, SE-41296, Gothenburg, Sweden

<sup>2</sup> Department of Mechanics and Maritime Sciences, Chalmers University of Technology, SE-41296, Gothenburg, Sweden

<sup>3</sup> Current affiliation: Volvo Car Corporation AB, Safety Center, SE-40531, Gothenburg, Sweden

\*chetsung@chalmers.se;alexofntu@gmail.com

# Outline

- **Motivation**
  - **Introduction**
  - **Do CNNs learn from texture or shape?**
  - **Methodology**
  - **Experimental Results**
  - **Conclusion and Future work**
-

# Motivation

- According to the National Highway Traffic Safety Administration, 3142 people were killed in motor vehicle crashes involving distracted drivers in 2019.
- Research on distraction has largely benefit from the analysis of Naturalistic Driving Data (NDD) in the last 15 years.
- One of the main concerns associated to the use of NDD is the time intensive and costly process of video reduction to extract variables from the videos.
- Convolutional Neural Networks (CNNs) are commonly thought to recognize objects by learning increasingly complex representations of object shapes. However, some recent studies suggested a more important role of image textures instead.
- Overfitting and over-confidence are two major issues that easily arise when training CNNs.

# Modeling driver monitoring as image classification

- Pros
  - 10 classes.
  - Most images were collected in real-driving scenario.
  - Collected in different cars, by people from different countries.
- Cons:
  - Captured with only two viewing angles

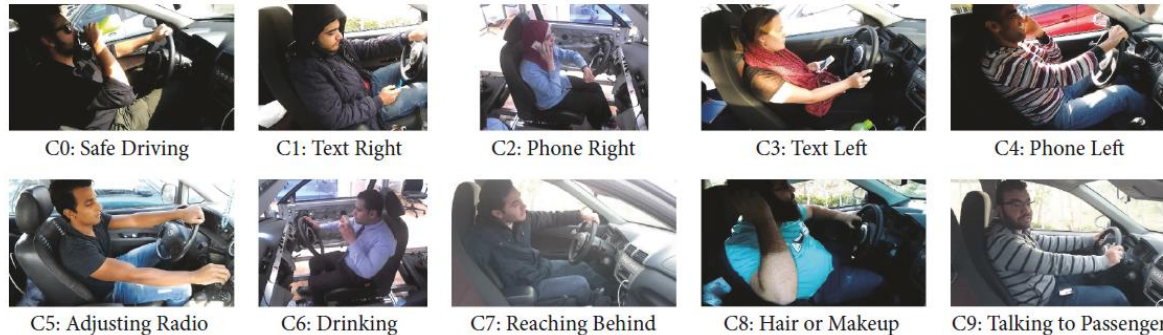
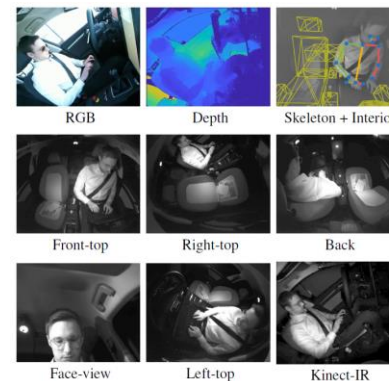


FIGURE 1: Ten classes of driver postures from our dataset.

# Drive & Act dataset

- Different views lead to different results and fully-seen view could provide better results.
- The more the driver's body captured by the camera; the higher the recognition accuracy.
- Model ensemble is beneficial.
- Pros:
  - The largest dataset of driver's action recognition
  - 6 different views
  - 3 modalities(RGB, IR and depth)
  - 83 classes
- Cons:
  - The same vehicle (Audi A3 ) for every driver in a simulated environment.
  - The number of drivers is limited. (10 for training, 2 for validation, 3 for testing)



Camera	View	Validation	Test
NIR Cameras	front top	69.57	63.64
	right top	65.16	60.80
	back	54.70	54.34
	face view	49.73	42.98
	left top	68.72	62.83
	combined	<u>72.70</u>	<u>67.17</u>
Kinect Color	right top	69.50	62.95
Kinect Depth		69.43	60.52
Kinect IR		72.90	64.98
Combined		<u>73.80</u>	<u>68.51</u>
All combined (score averaging)		<b><u>74.85</u></b>	<b><u>69.03</u></b>

# EuroFOT and Drive C2X

- Pros:
  - Provides significant more images with different views
  - 8 classes of distractions and 5 phone usages and four hands-on-wheel classes.
  - More drivers (127)
- Cons:
  - Monochrome images.
  - Limited resolution (352x288)

<b>Class and images</b>	<b>Training</b>	<b>Validation</b>
No activities	169213	200432
Interaction with passenger	1244	504
Talking or singing	18106	732
Reaching for an object	19232	8046
Interaction with center stack	6836	3439
Eating/Drinking	4975	2784
Hands-face interaction	29847	19598
Reading	1920	1058

# CNN learns from texture more!

- This paper has shown that ImageNet-trained CNNs are strongly biased towards recognizing textures rather than shapes.

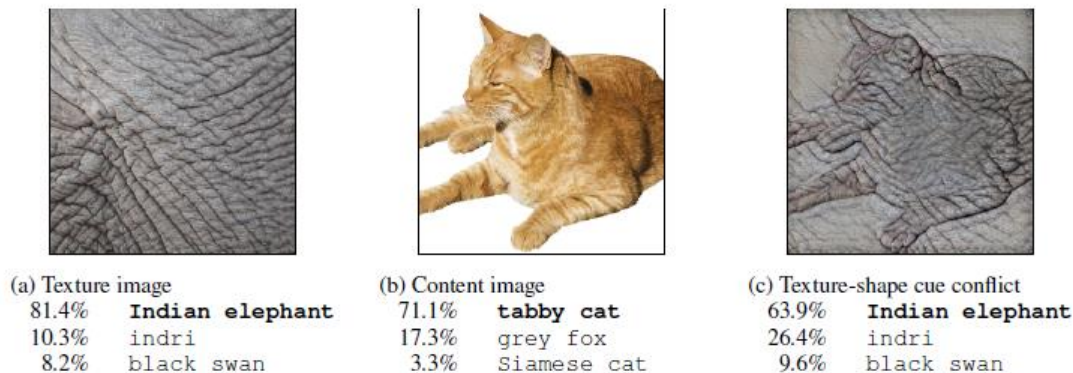


Figure 1: Classification of a standard ResNet-50 of (a) a texture image (elephant skin: only texture cues); (b) a normal image of a cat (with both shape and texture cues), and (c) an image with a texture-shape cue conflict, generated by style transfer between the first two images.



# CNN's performance degrades when lacking of texture information

- Pretraining with stylized-Imagenet is beneficial in terms of image classification and object detection since a shape-based representation is more beneficial than a texture-based representation.
- However, due to copyright issue, stylized-Imagenet can't be directly provided.
- The neural-style transfer model needs 1 min (500 iterations) for each image.
- For 1.2 million images, we need 833 days.



Stylized-ImageNet

name	training	fine-tuning	top-1 IN accuracy (%)	top-5 IN accuracy (%)	Pascal VOC mAP50 (%)
vanilla ResNet	IN	-	76.13	92.86	70.7
	SIN	-	60.18	82.62	70.6
	SIN+IN	-	74.59	92.14	74.0
Shape-ResNet	SIN+IN	IN	76.72	93.28	75.1

Imagenet classification and VOC detection results

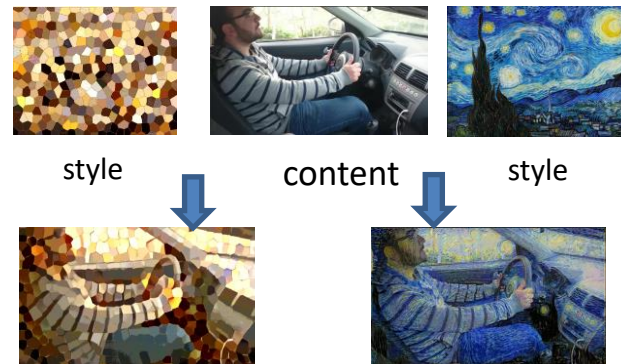


## Neural-style transfer as a data augmentation strategy for DDD

- The number of Images in Distracted Driver Dataset: 14478 images.
- By reducing iterations to 200 iterations, we need 3.35 days.



FIGURE 1: Ten classes of driver postures from our dataset.



## Label smoothing in image classification

- Label Smoothing (Rafael et al., 2019) is a regularization technique that introduces noise for the labels. This accounts for the fact that datasets may have some mistakes in them, so maximizing the likelihood of  $\log p(y|x)$  might result in over-fitting.
- Assume for a small constant  $\epsilon$ , the target value of the training label  $y$  is  $1-\epsilon$  and  $\epsilon/(k-1)$  for the target class and others, respectively. i.e., the original target value of each class is

$$P_i = \begin{cases} 1, & i = y, \\ 0, & i \neq y. \end{cases} \quad (1)$$

After label-smoothing, they become

$$P_i = \begin{cases} 1 - \epsilon, & i = y, \\ \epsilon/(k - 1), & i \neq y. \end{cases} \quad (2)$$

Therefore, for cross-entropy loss

$$Loss = -\sum_{i=1}^k p_i \log q_i, \quad (3)$$

## Experimental Results

- Under the same backbone-ResNet-50, Label Smoothing is always helpful if cross-entropy is applied.
- The backbone trained by Stylized-ImageNet, learns to capture shapes instead texture so that the model is more robust against noise.
- The backbone also generalizes to EuroFOT well.

*Table 3 ResNet-50 results on Distracted Driver Dataset*

Pretraining	Finetuning	Label smoothing	Accuracy
ImageNet	DDD	No	81.695% (Eraqi et al., 2019)
ImageNet	DDD	Yes	86.95%
ImageNet + Stylized-ImageNet	DDD	Yes	88.05%
ImageNet + Stylized-ImageNet	DDD+ Stylized-DDD	Yes	89.01%
ImageNet + Stylized-ImageNet	EuroFOT	Yes	84.72%

## Conclusion and Future work

- We expect to apply the same model on phone-usage and hands-on-wheel classification.
- Most driving monitoring application is actually a multi-label image classification application, i.e., the driver might actually be talking (a distraction behavior) and using smartphone with his left hand (a phone usage class) simultaneously.
- Large-scale ensemble is actually an option when this application is applied on a labelling system.
- Driver distraction classification can actually be trained in semi-supervised setting via mean teacher model.



**CHALMERS**  
UNIVERSITY OF TECHNOLOGY

Thank you for listening!