

Robotic Vision

E.M. Bakker

Overview

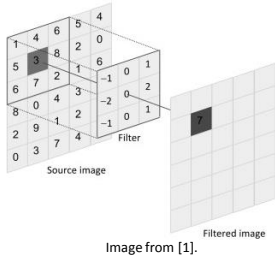
- OpenCV
- Some Neural Networks and AlexNet

Computer Vision and Pattern Recognition (CVPR)

- Object Tracking
- Human Robot Interaction
- Some problems with Neural Networks
- ...

OpenCV

- Low level image processing.
- Convolutional Kernels: filters, edge detectors, etc.



The general expression of a convolution is

$$g(x, y) = \omega * f(x, y) = \sum_{s=-a}^a \sum_{t=-b}^b \omega(s, t) f(x-s, y-t),$$

where $g(x, y)$ is the filtered image, $f(x, y)$ is the original image, ω is the filter kernel. Every element of the filter kernel is considered by $-a \leq s \leq a$ and $-b \leq t \leq b$.

Wikipedia

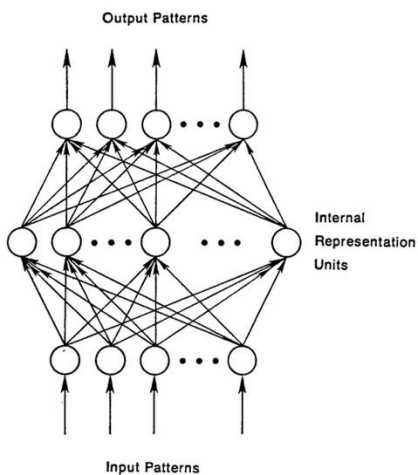
- Blob tracking
- Face and people detector
- Neural networks

[1] <https://www.sciencedirect.com/topics/computer-science/convolution-filter>

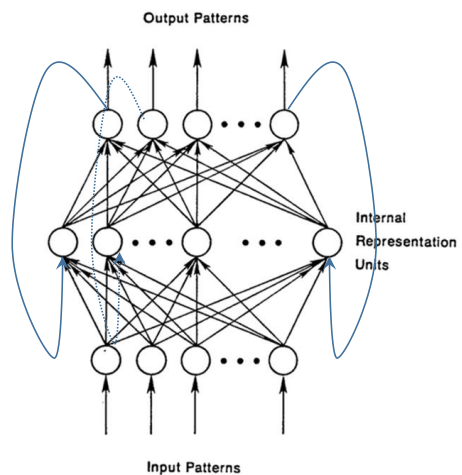
Operation	Kernel ω	Image result $g(x,y)$
Identity	$\begin{bmatrix} 0 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{bmatrix}$	
Edge detection	$\begin{bmatrix} 1 & 0 & -1 \\ 0 & 0 & 0 \\ -1 & 0 & 1 \end{bmatrix}$	
	$\begin{bmatrix} 0 & 1 & 0 \\ 1 & -4 & 1 \\ 0 & 1 & 0 \end{bmatrix}$	
	$\begin{bmatrix} -1 & -1 & -1 \\ -1 & 8 & -1 \\ -1 & -1 & -1 \end{bmatrix}$	
Sharpen	$\begin{bmatrix} 0 & -1 & 0 \\ -1 & 5 & -1 \\ 0 & -1 & 0 \end{bmatrix}$	
Box blur (normalized)	$\frac{1}{9} \begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{bmatrix}$	

Wikipedia

Some Neural Networks



Feed Forward Neural Network



Recurrent Neural Network

DNN: AlexNet, VGG16, ResNet, etc.

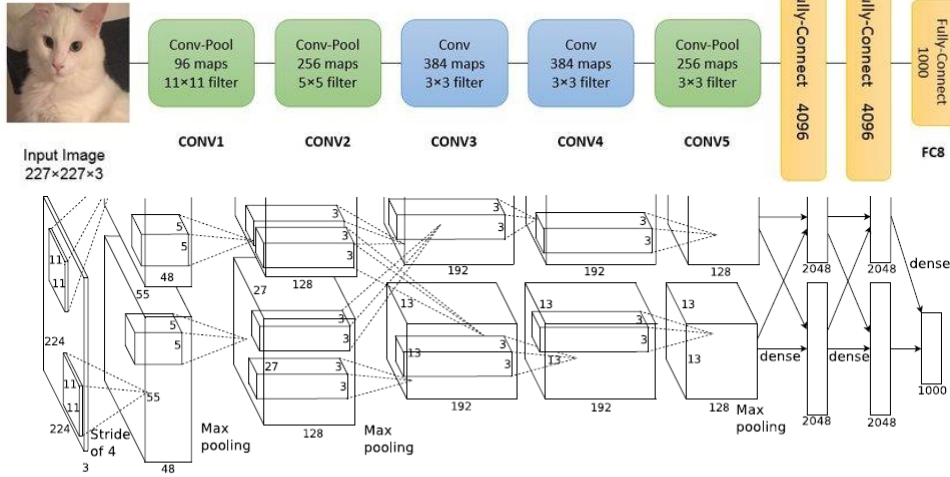


Figure 2: An illustration of the architecture of our CNN, explicitly showing the delineation of responsibilities between the two GPUs. One GPU runs the layer-parts at the top of the figure while the other runs the layer-parts at the bottom. The GPUs communicate only at certain layers. The network's input is 150,528-dimensional, and the number of neurons in the network's remaining layers is given by 253,440–186,624–64,896–64,896–43,264–4096–4096–1000.

Krizhevsky, Alex; Sutskever, Ilya; Hinton, Geoffrey E. "ImageNet classification with deep convolutional neural networks" Communications of the ACM. 60 (6): 84–90.

Deep Visualization Toolbox

yosinski.com/deepvis

#deepvis



Jason Yosinski



Jeff Clune



Anh Nguyen



Thomas Fuchs



Hod Lipson



Object Tracking

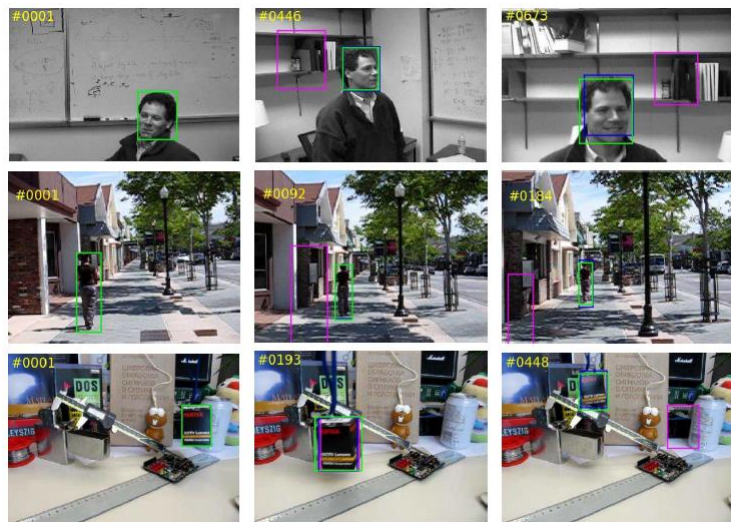
- Conference on Computer Vision and Pattern Recognition (CVPR)

Real-Time Tracking

- A. He et al. A Twofold Siamese Network for Real-Time Object Tracking
- B. Yang et al. PIXOR: Real-Time 3D Object Detection From Point Clouds
- B. Tekin et al., Real-Time Seamless Single Shot 6D Object Pose Prediction

A. He et al. A Twofold Siamese Network for Real-Time Object Tracking, CVPR2018.

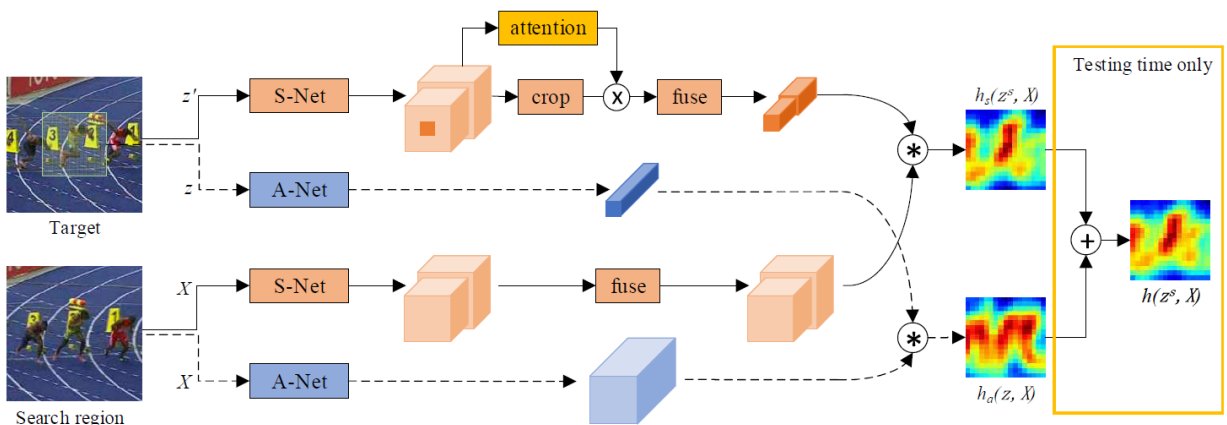
- **Green** is ground truth.
- **Purple** is tracked by *SiamFC*.
- **Blue** is tracked by the novel twofold Siamese network *2FSiamFC*.
- *2FSiamFC* is more robust to shooting angle change and scale change.



A. He et al. A Twofold Siamese Network for Real-Time Object Tracking, CVPR2018.

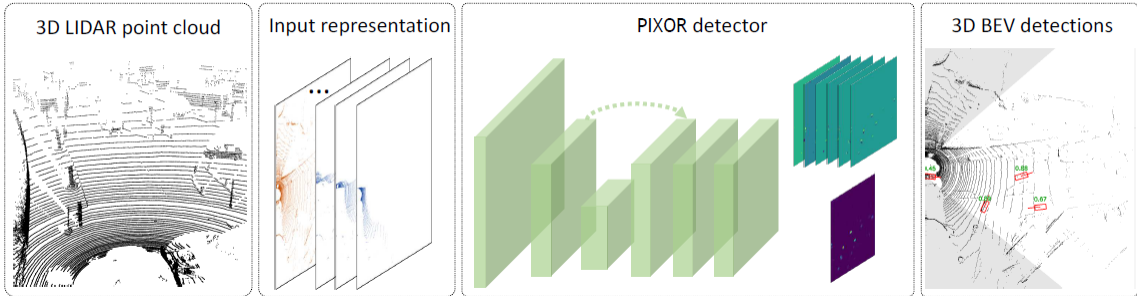
Object Tracking is a **similarity learning problem**

- Compare the target image patch with the candidate patches in a search region.
- Track the object to the location where the highest similarity score is obtained.
- Similarity learning with deep CNNs is done using so called Siamese architectures (SiamFC).
- CNNs can process a larger search image where all sub-windows are evaluated as similarity candidates. (Efficient.)



- A-Net is an appearance network.
- The dotted lines is a SiamFC (Fully Convolutional Siamese Network Bertinetto et al. 2016.)
- S-Net is a semantic network.
- The channel attention module determines the weight for each feature channel based on both target and context information.
- Appearance branch and Semantic branch are trained separately.

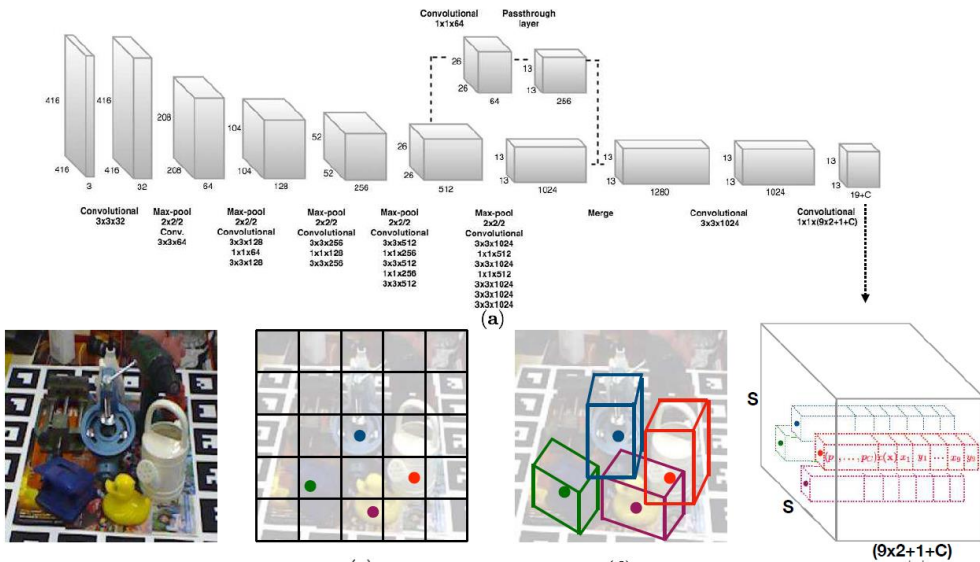
B. Yang et al. PIXOR: Real-Time 3D Object Detection From Point Clouds (CVPR2018)



PIXOR, a proposal-free, single-stage detector that outputs oriented 3D object estimates decoded from pixel-wise neural network predictions.

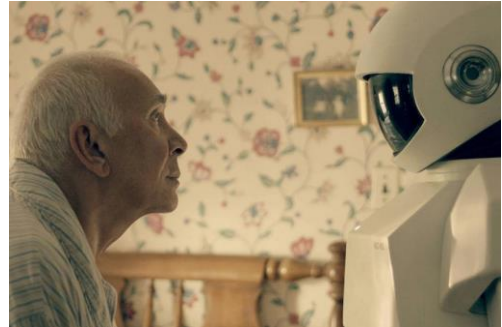
- real-time 3D object detection from point clouds in the context of autonomous driving.
- 3D data by representing the scene from the *Bird's Eye View (BEV)*
- **Evaluation 10fps state-of-the-art:** using the KITTI BEV object detection benchmark, and a large-scale 3D vehicle detection benchmark.

Tekin et al., Real-Time Seamless Single Shot 6D Object Pose Prediction (CVPR2018)



Human Robot Interaction

- Face Recognition
- Pose Recognition
- Hand Tracking
- Person Tracking
- Emotion Recognition
- Action Recognition



Face Recognition

- Yancheng Bai, et al., Finding Tiny Faces in the Wild With Generative Adversarial Network, CVPR, 2018.
- Xuanyi Dong, et al., Aggregated Network for Facial Landmark Detection, CVPR, 2018.
- Yaojie Liu, et al., Learning Deep Models for Face Anti-Spoofing: Binary or Auxiliary Supervision, CVPR, 2018.
- CVPR2018 58 papers on Face Recognition
- CVPR2019 similar

Yancheng Bai, et al., Finding Tiny Faces in the Wild With Generative Adversarial Network, CVPR, 2018.

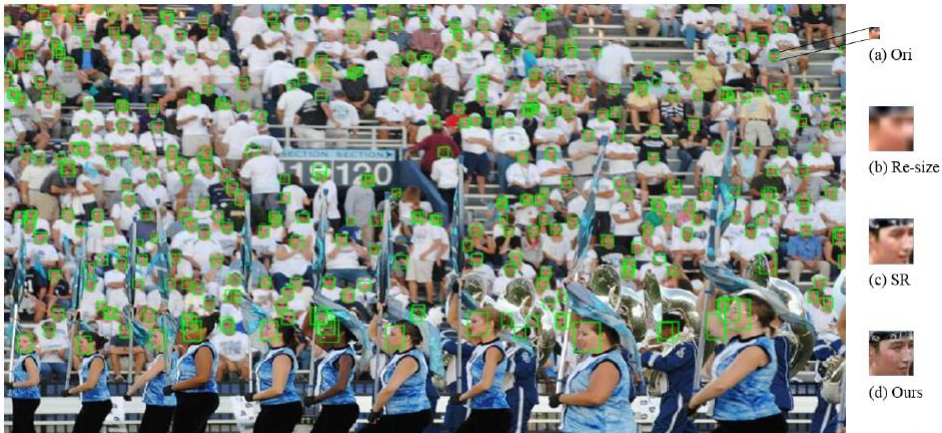
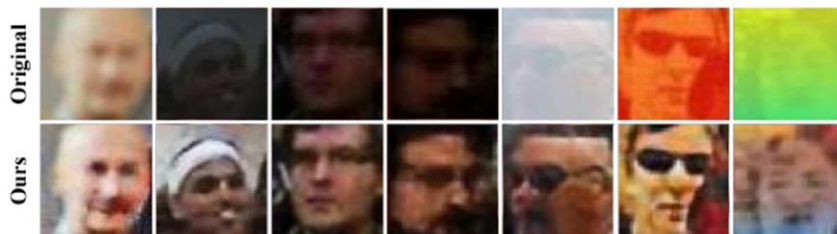
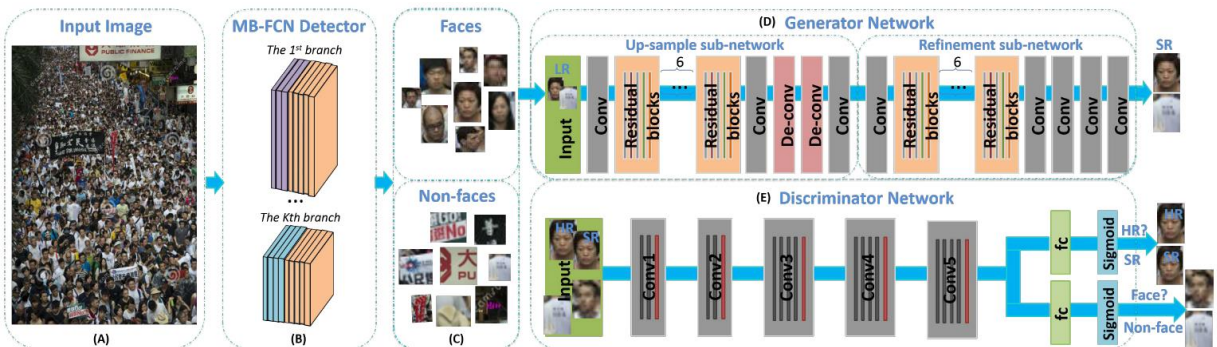


Figure1. The detection results of tiny faces in the wild. (a) is the original low-resolution blurry face, (b) is the result of re-sizing directly by a bi-linear kernel, (c) is the generated image by the super-resolution method, and our result (d) is learned by the super-resolution ($\times 4$ upscaling) and refinement network simultaneously. Best viewed in color and zoomed in.

Generative Adversarial Network.



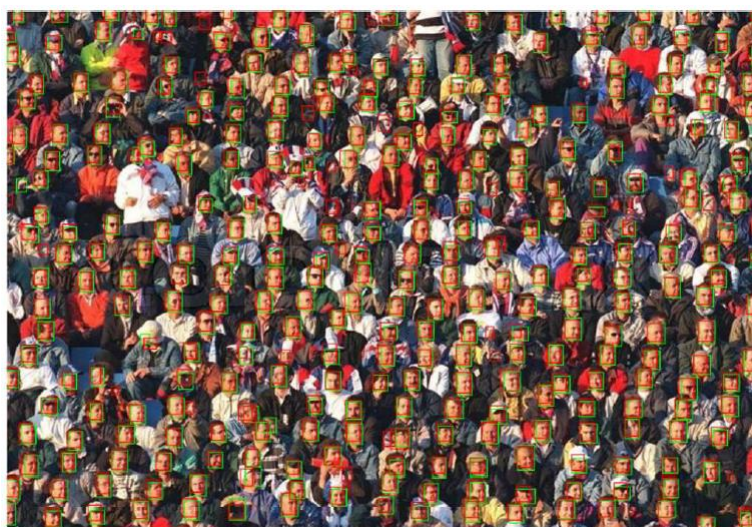
Some Qualitative Results

Green ground truth, red selected by the network.



Some Qualitative Results

Green ground truth, red selected by the network.



Some Qualitative Results

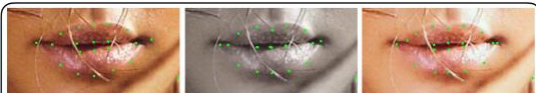
Green ground truth, red selected by the network.



Xuanyi Dong, et al., Aggregated Network for Facial Landmark Detection, CVPR2018.



landmark predictions on the *same* image
under *different* styles

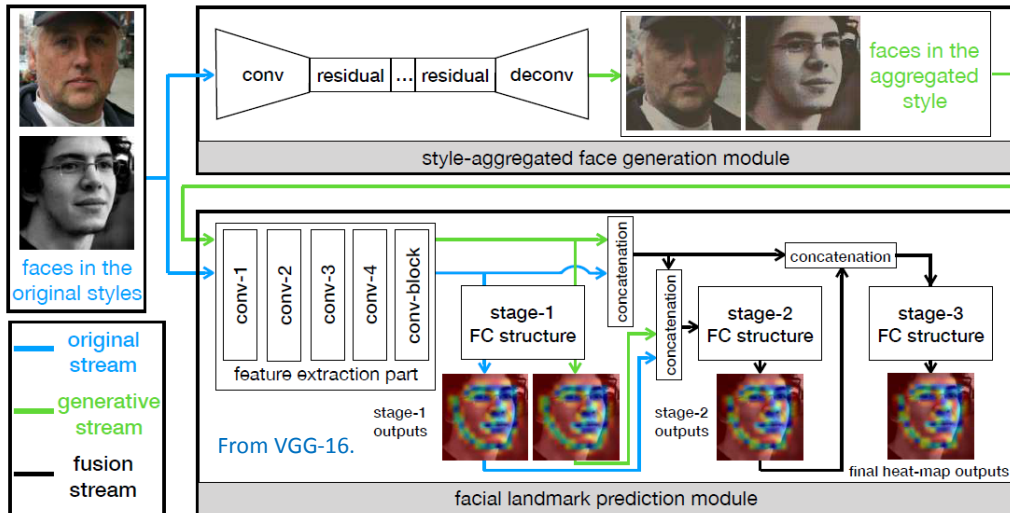


landmark predictions on zoom-in parts



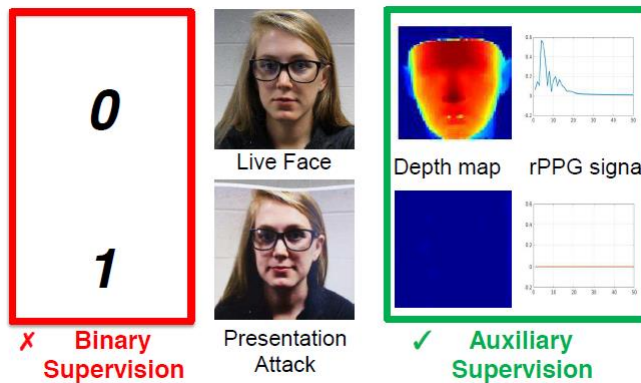
- Landmark detection robust to 'style' changes: dark, light, color, gray. (300-W data set.)

Xuanyi Dong, et al., Aggregated Network for Facial Landmark Detection, CVPR, 2018.



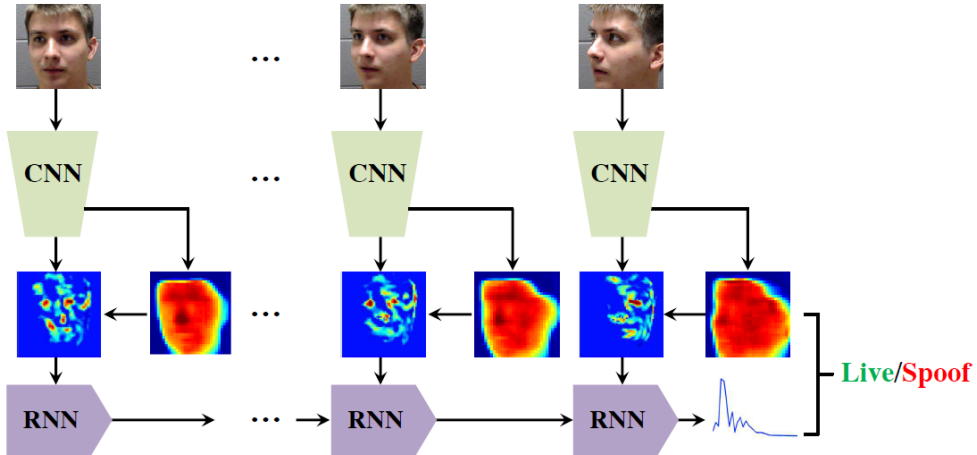
Landmark detection robust to 'style' changes: dark, light, color, gray. (300-W data set.) Network uses a GAN for generated aggregated styles (green lines).

Yaojie Liu, et al., Learning Deep Models for Face Anti-Spoofing: Binary or Auxiliary Supervision, CVPR2018.



Use extra, live information for face recognition in order to make spoofing more difficult.

Yaojie Liu, et al., Learning Deep Models for Face Anti-Spoofing: Binary or Auxiliary Supervision, CVPR, 2018.



S. Zhang et al.
A Dataset and Benchmark
for Large-scale Multi-modal
Face Anti-spoofing.
CVPR2019

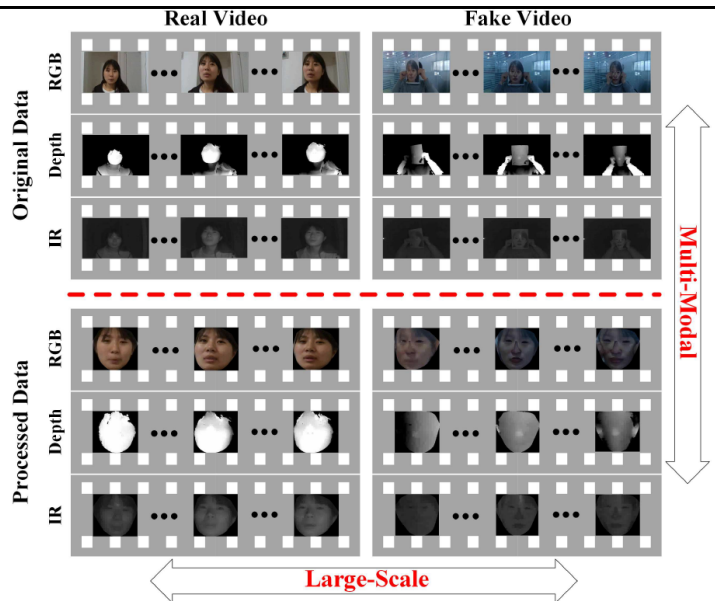


Figure 1. The CASIA-SURF dataset. It is a large-scale and multi-modal dataset for face anti-spoofing, consisting of 492,522 images with 3 modalities (*i.e.*, RGB, Depth and IR).

A. Tonioni et al. Real-time self-adaptive deep stereo. CVPR2019
<https://github.com/CVLAB-Unibo/Real-time-self-adaptive-deep-stereo>

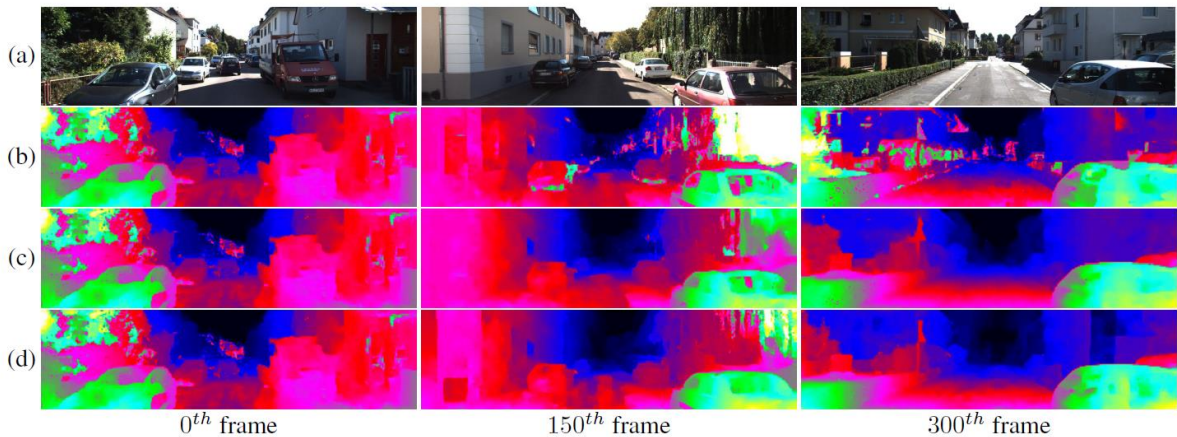


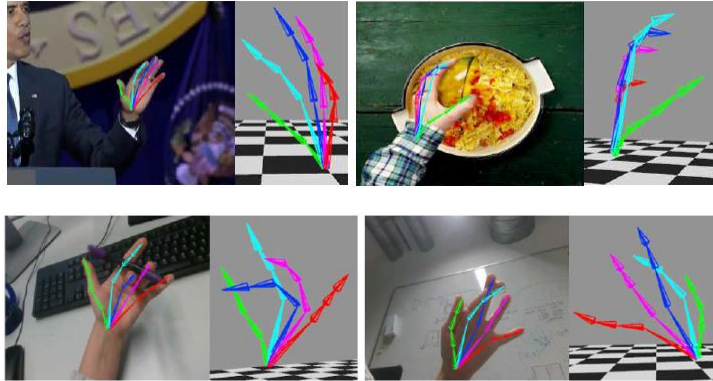
Figure 1. Disparity maps predicted by *MADNet* on a KITTI sequence [7]. Left images (a), no adaptation (b), online adaptation of the *whole* network (c), online adaptation by *MAD* (d). Green pixel values indicate larger disparities (*i.e.*, closer objects).

Hand Pose Recognition

F. Mueller, et al., **GANerated Hands for Real-Time 3D Hand Tracking From Monocular RGB**, CVPR2018.

G. Garcia-Hernando, et al., **First-Person Hand Action Benchmark With RGB-D Videos and 3D Hand Pose Annotations**, CVPR2018.

F. Mueller, et al., GANerated Hands for Real-Time 3D Hand Tracking From Monocular RGB, CVPR2018.

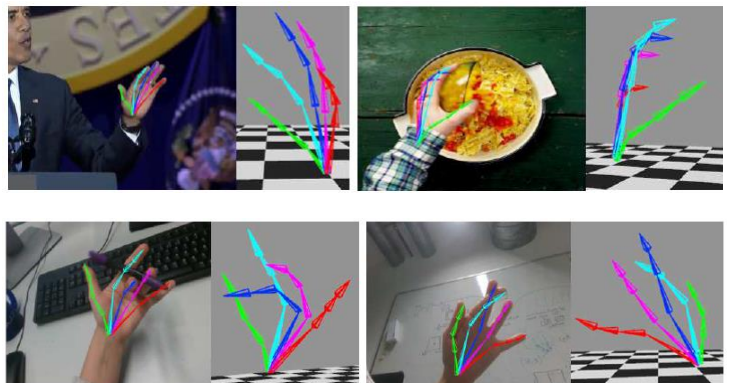


Input: RGB Image
Output: Hand Pose Skeleton.

F. Mueller, et al., GANerated Hands for Real-Time 3D Hand Tracking From Monocular RGB, CVPR, 2018.

Real-time 3D hand tracking from monocular RGB-only input.

- Works on unconstrained videos from YouTube
- Is robust to occlusions.
- Real-time 3D hand tracking using an off-the-shelf RGB webcam in unconstrained setups.



F. Mueller, et al., GANerated Hands for Real-Time 3D Hand Tracking From Monocular RGB, CVPR, 2018.

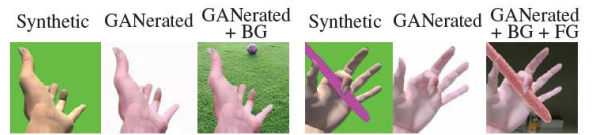


Figure 5: Two examples of synthetic images with background/object masks in green/pink.

- **GeoConGAN** produces ‘real’ images from synthetic images. These ‘real’ images are then used to train **RegNet**.
- The trained **RegNet** is used to recognize global 3d hand poses in real time from RGB video streams.

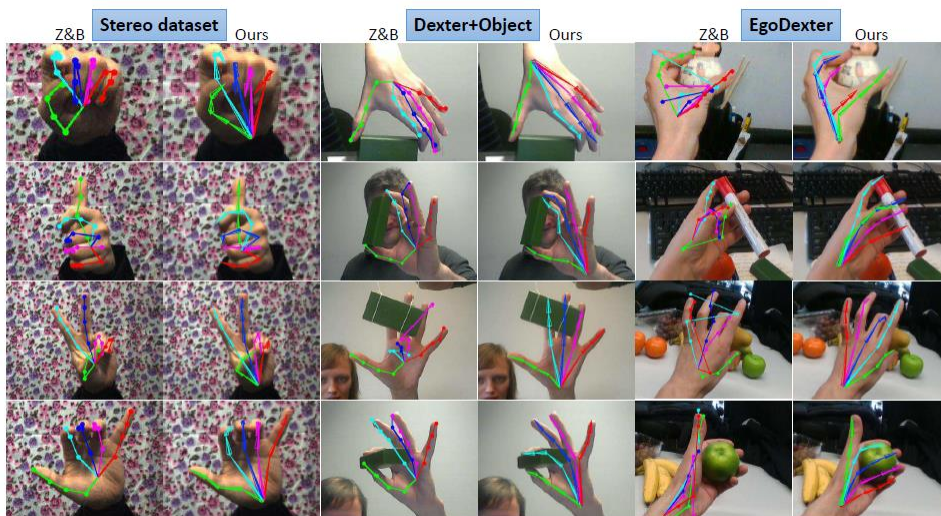
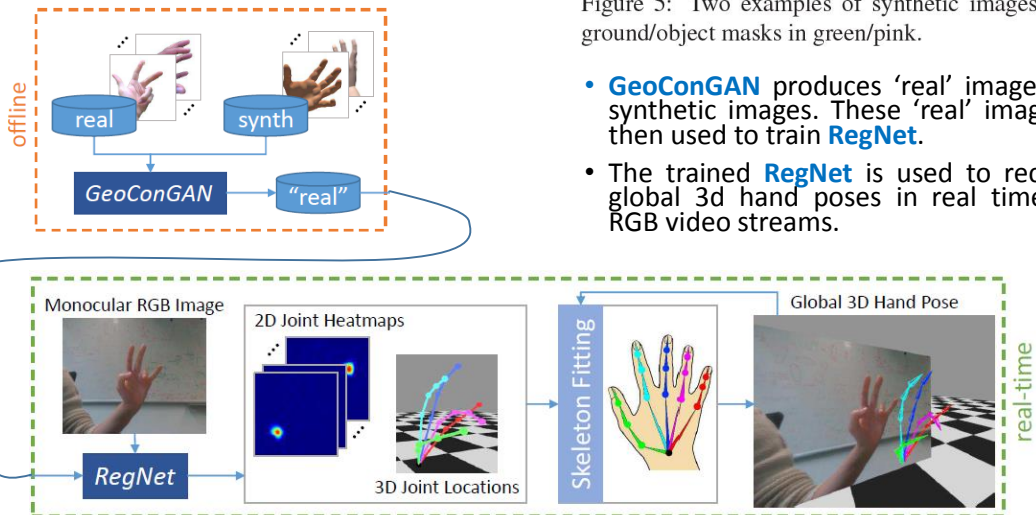
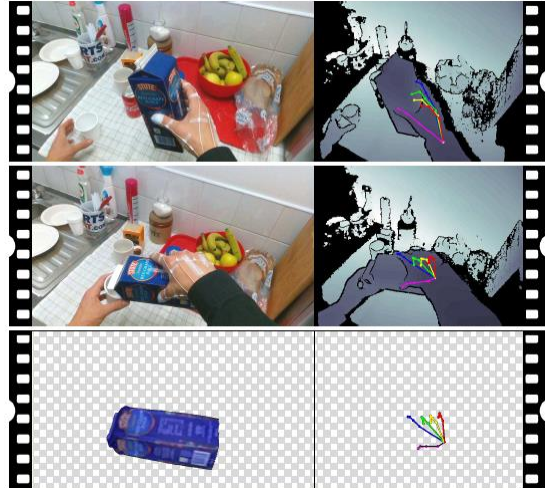


Figure 8: We compare our results with Zimmermann and Brox [63] on three different datasets. Our method is more robust in cluttered scenes and it even correctly retrieves the hand articulation when fingers are hidden behind objects.

Garcia-Hernando, et al., **First-Person Hand Action Benchmark With RGB-D Videos and 3D Hand Pose Annotations**, CVPR2018.

Pouring Juice

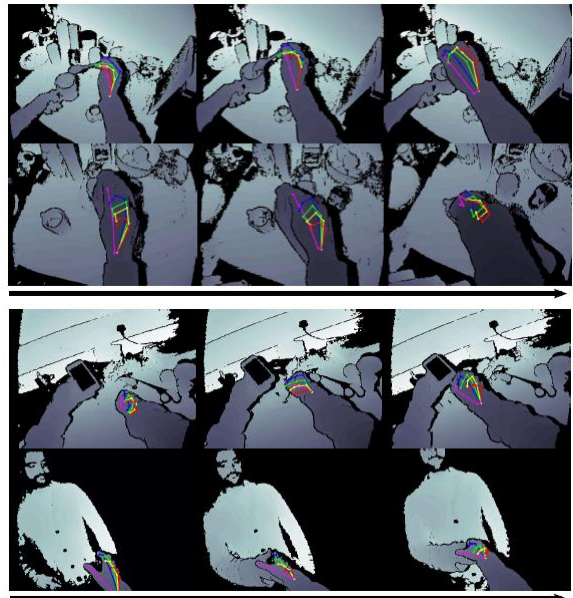
- A novel firstperson action recognition dataset with RGB-D videos and 3D hand pose annotations.
- Magnetic sensors and inverse kinematics to capture the hand pose.
- Also captured 6D object pose for some of the actions



Garcia-Hernando, et al., **First-Person Hand Action Benchmark With RGB-D Videos and 3D Hand Pose Annotations**, CVPR, 2018.

A novel first person action recognition dataset with RGB-D videos and 3D hand pose annotations.

- Put sugar.
- Pour milk.
- Charge cell-phone.
- Shake hand.



Garcia-Hernando, et al., First-Person Hand Action Benchmark With RGB-D Videos and 3D Hand Pose Annotations, CVPR, 2018.

Visual data: Intel RealSense SR300 RGB-D camera on the shoulder of the subject (RGB 30 fps at 1920×1080 and Depth 640×480.)

Pose annotation:

hand pose

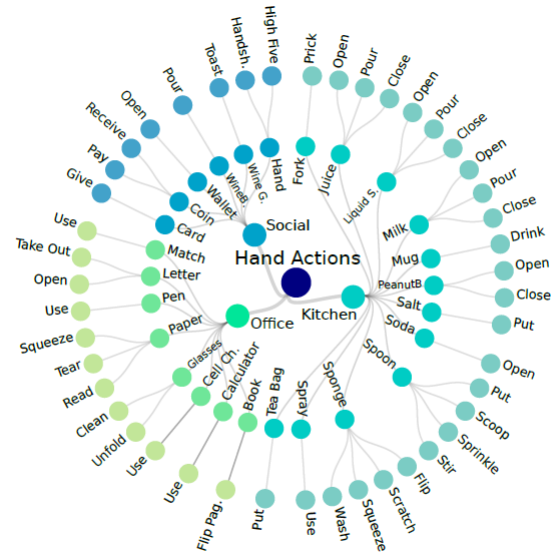
- captured using six magnetic sensors (6DOF) attached to the user's hand, five fingertips and one wrist, following [84].
- the hand pose is inferred using inverse kinematics over a defined 21-joint hand model

object pose

- 1 6DOF magnetic sensor attached to the closest point to the center of mass.

Recording process:

- 6 people, all right handed performed the actions.



Garcia-Hernando, et al., First-Person Hand Action Benchmark With RGB-D Videos and 3D Hand Pose Annotations, CVPR2018.

Baseline: RNN LSTM 100 neurons.

1:3 25% training 75% testing

1:1 50% - 50%

3:1 75% - 25%

Cross-person

Leave one of the 6 persons out of the training and test on the person left out.

Tensorflow and Adam optimizer.

Baseline Action recognition results

Protocol	1:3	1:1	3:1	cross-person
Acc. (%)	58.75	78.73	84.82	62.06

Hand pose recognition

Method	Year	Color	Depth	Pose	Acc. (%)
Two stream-color [15]	2016	✓	✗	✗	61.56
Two stream-flow [15]	2016	✓	✗	✗	69.91
Two stream-all [15]	2016	✓	✗	✗	75.30
HOG ² -depth [40]	2013	✗	✓	✗	59.83
HOG ² -depth+pose [40]	2013	✗	✓	✓	66.78
HON4D [43]	2013	✗	✓	✗	70.61
Novel View [47]	2016	✗	✓	✗	69.21
1-layer LSTM	2016	✗	✗	✓	78.73
2-layer LSTM	2016	✗	✗	✓	80.14
Moving Pose [85]	2013	✗	✗	✓	56.34
Lie Group [64]	2014	✗	✗	✓	82.69
HBRNN [12]	2015	✗	✗	✓	77.40
Gram Matrix [86]	2016	✗	✗	✓	85.39
TF [17]	2017	✗	✗	✓	80.69
JOULE-color [19]	2015	✓	✗	✗	66.78
JOULE-depth [19]	2015	✗	✓	✗	60.17
JOULE-pose [19]	2015	✗	✗	✓	74.60
JOULE-all [19]	2015	✓	✓	✓	78.78

Table 4: Hand action recognition performance by different evaluated approaches on our proposed dataset.

Some Problems with Deep Neural Networks

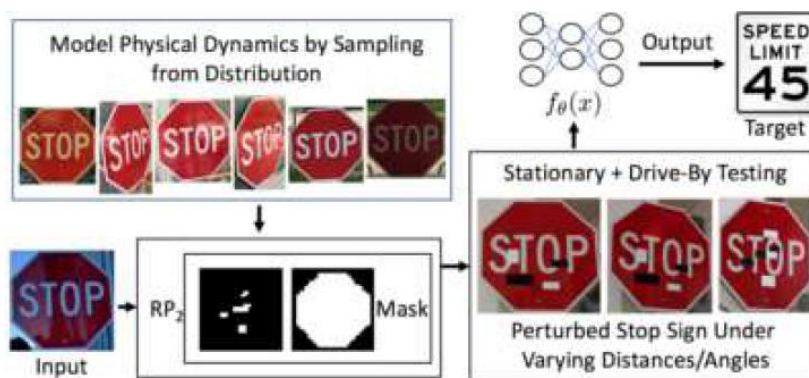
K. Eykholt, et al. Dawn Song Robust Physical-World Attacks on Deep Learning Visual Classification, CVPR2018.



K. Eykholt, et al. Dawn Song Robust Physical-World Attacks on Deep Learning Visual Classification, CVPR2018.

Robust Physical Perturbations (RP2):

- generate physical perturbations for physical-world objects such that a DNN-based classifier produces a designated misclassification.
- This under a range of dynamic physical conditions, including different viewpoint angles and distances.



K. Eykholt, et al. Dawn Song Robust Physical-World Attacks on Deep Learning Visual Classification, CVPR2018.

Two types of attacks showing that RP2 produces robust perturbations for real road signs.

- poster attacks are successful in 100% of stationary and drive-by tests against LISA-CNN
- sticker attacks are successful in 80% of stationary testing conditions



K. Eykholt, et al. Dawn Song Robust Physical-World Attacks on Deep Learning Visual Classification, CVPR2018.



This is a micro-wave.

This is not a micro-wave.



References

Papers can be obtained from <http://openaccess.thecvf.com/CVPR2018.py>

Real-Time Tracking

- [1] A. He et al. A Twofold Siamese Network for Real-Time Object Tracking, CVPR, 2018.
- [2] B. Yang et al. PIXOR: Real-Time 3D Object Detection From Point Clouds, CVPR, 2018.
- [3] B. Tekin et al., Real-Time Seamless Single Shot 6D Object Pose Prediction, CVPR, 2018.

Face Recognition

- [4] Yancheng Bai, et al., Finding Tiny Faces in the Wild With Generative Adversarial Network, CVPR, 2018.
- [5] Xuanyi Dong, et al., Aggregated Network for Facial Landmark Detection, CVPR, 2018.
- [6] Yaojie Liu, et al., Learning Deep Models for Face Anti-Spoofing: Binary or Auxiliary Supervision, CVPR, 2018.

Hand Pose Recognition

- [7] F. Mueller, et al., GANerated Hands for Real-Time 3D Hand Tracking From Monocular RGB, CVPR, 2018.
- [8] G. Garcia-Hernando, et al., First-Person Hand Action Benchmark With RGB-D Videos and 3D Hand Pose Annotations, CVPR, 2018.

Problems with Deep Learning Classification

- [9] K. Eykholt, et al. Dawn Song Robust Physical-World Attacks on Deep Learning Visual Classification, CVPR, 2018.

Conference on Computer Vision and Pattern Recognition

- <http://openaccess.thecvf.com/CVPR2018.py>
- <http://openaccess.thecvf.com/CVPR2019.py>