# Multi-sensor System for Driver's Hand-Gesture Recognition

Pavlo Molchanov, Shalini Gupta,
Kihwan Kim, and Kari Pulli

# Driver distraction

Current interfaces in cars distract drivers from the road

# Driver distraction



(http://www.softkinetic.com)

Touchless interfaces will help keep drivers' attention on the road
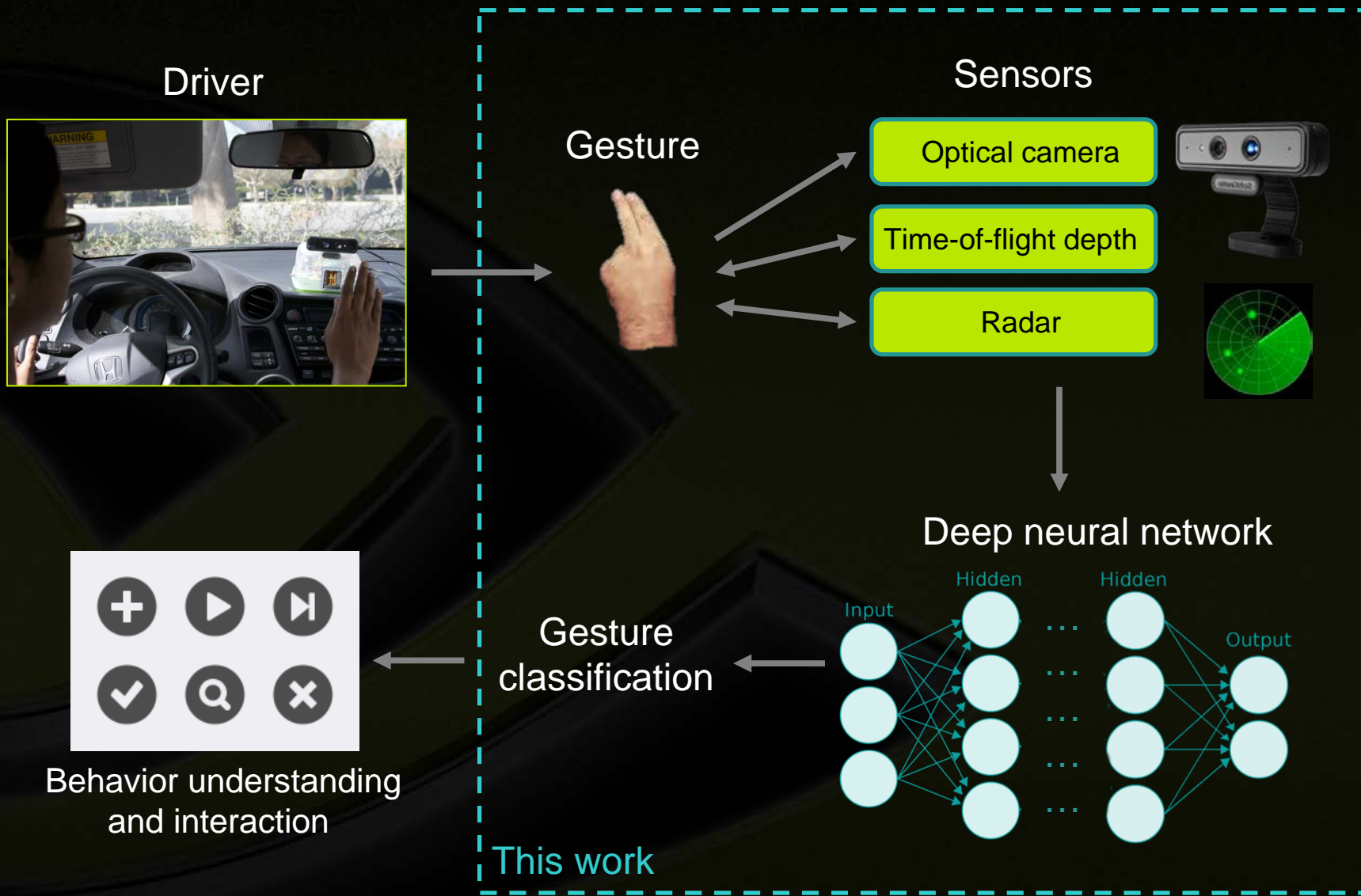
# Existing work

- Neverova et al.*, ECCV ChaLearn Workshop* 2014.
    - RGBD and upper body skeletal pose
    - Deep neural networks
    - Indoors only

- Ohn-Bar and Trivedi, *IEEE Trans. ITS* 2014.
    - RGBD
    - HOG+HOG$^2$ and SVM classifier
    - In car during day and evening

# Our solution

Driver

Sensors

Gesture

Optical camera

Time-of-flight depth

Radar

Deep neural network

Input
Hidden
Hidden
Output

Gesture classification

Behavior understanding and interaction

This work

5

# Why multi-sensor?

**Day**

**Night**

Color cameras

# Why multi-sensor?

Color cameras

**Day**

**Night**

Commodity depth
cameras

**No sunlight**

**Sunlight**

7

# Radar sensor

All lighting conditions:

No interference from the Sun:

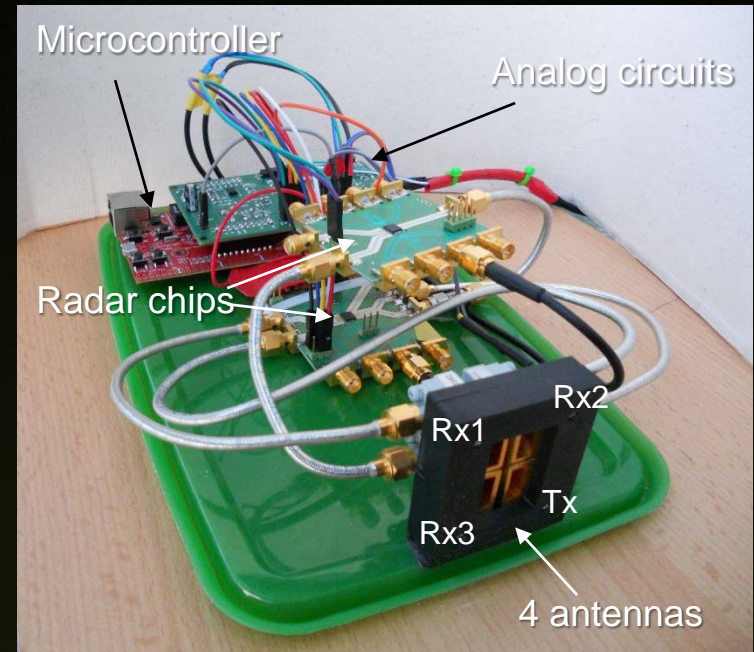Direct measurements of local radial velocities (by Doppler shift):
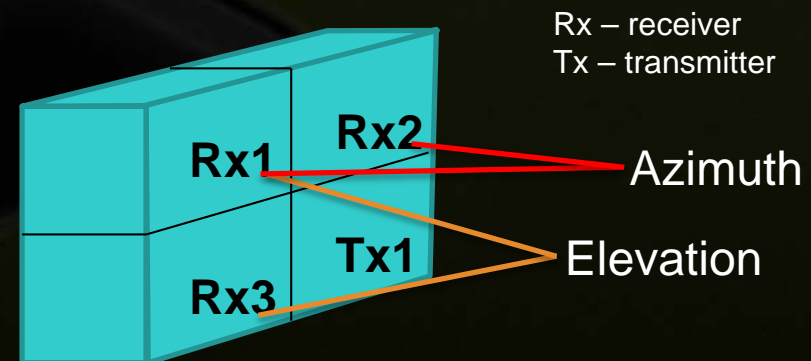
$v$

# Radar sensor

- Frequency Modulated Continuous wave (FMCW) radar architecture, 24 GHz

- Our design

- Molchanov et al., *IEEE Radar*, 2015.

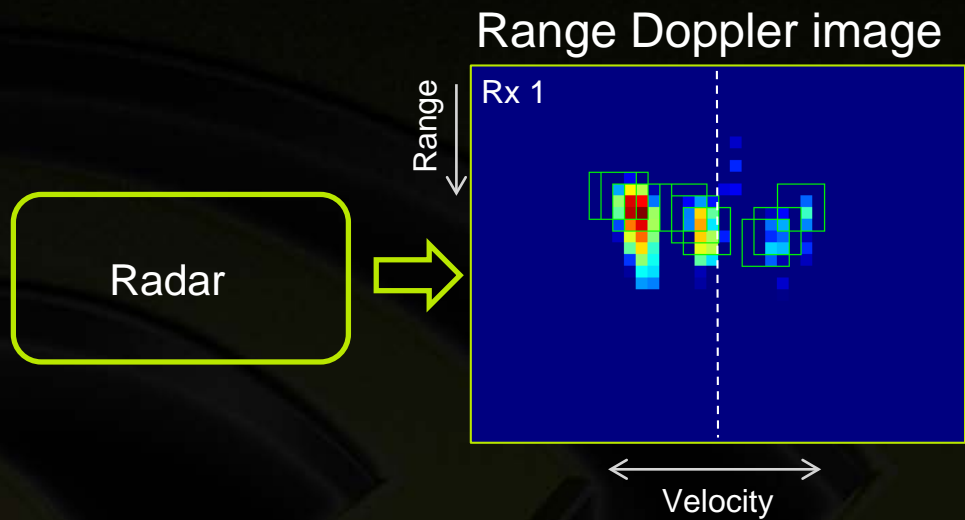## Radar can estimate:

- Range with resolution of <u>4 cm</u>

- Radial velocity with resolution <u>0.04 m/s</u>

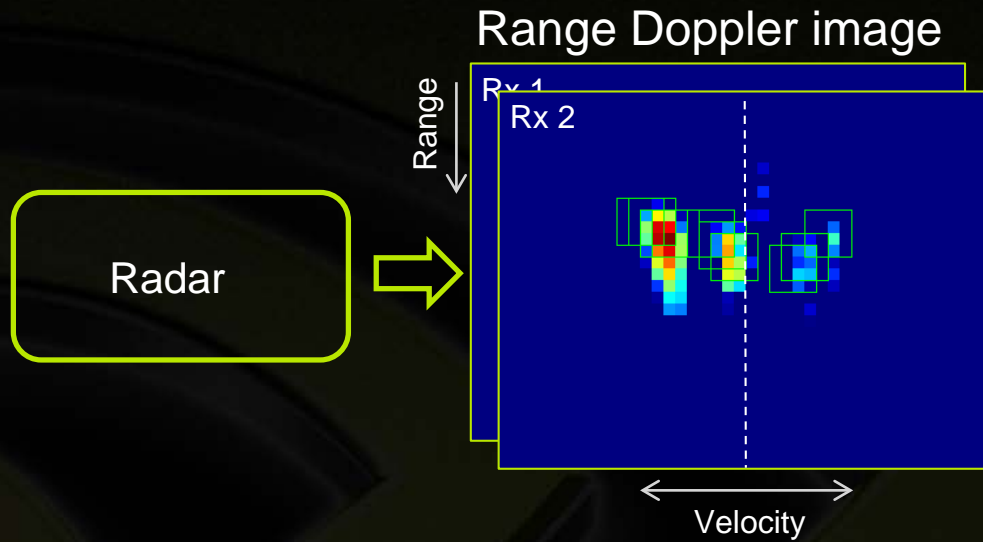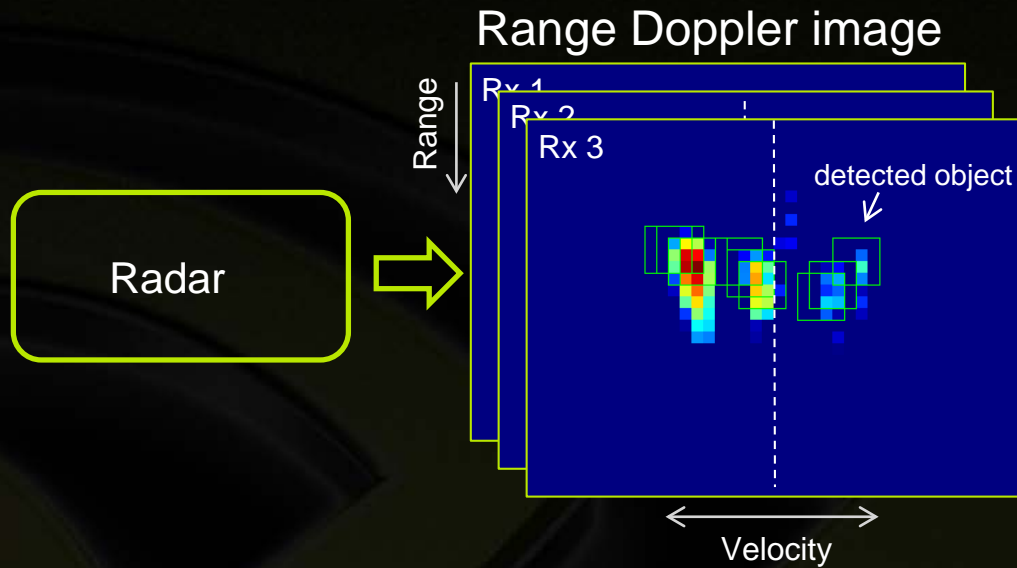- Angles of arriving (<u>azimuth</u> and <u>elevation</u>) are estimated for detected objects
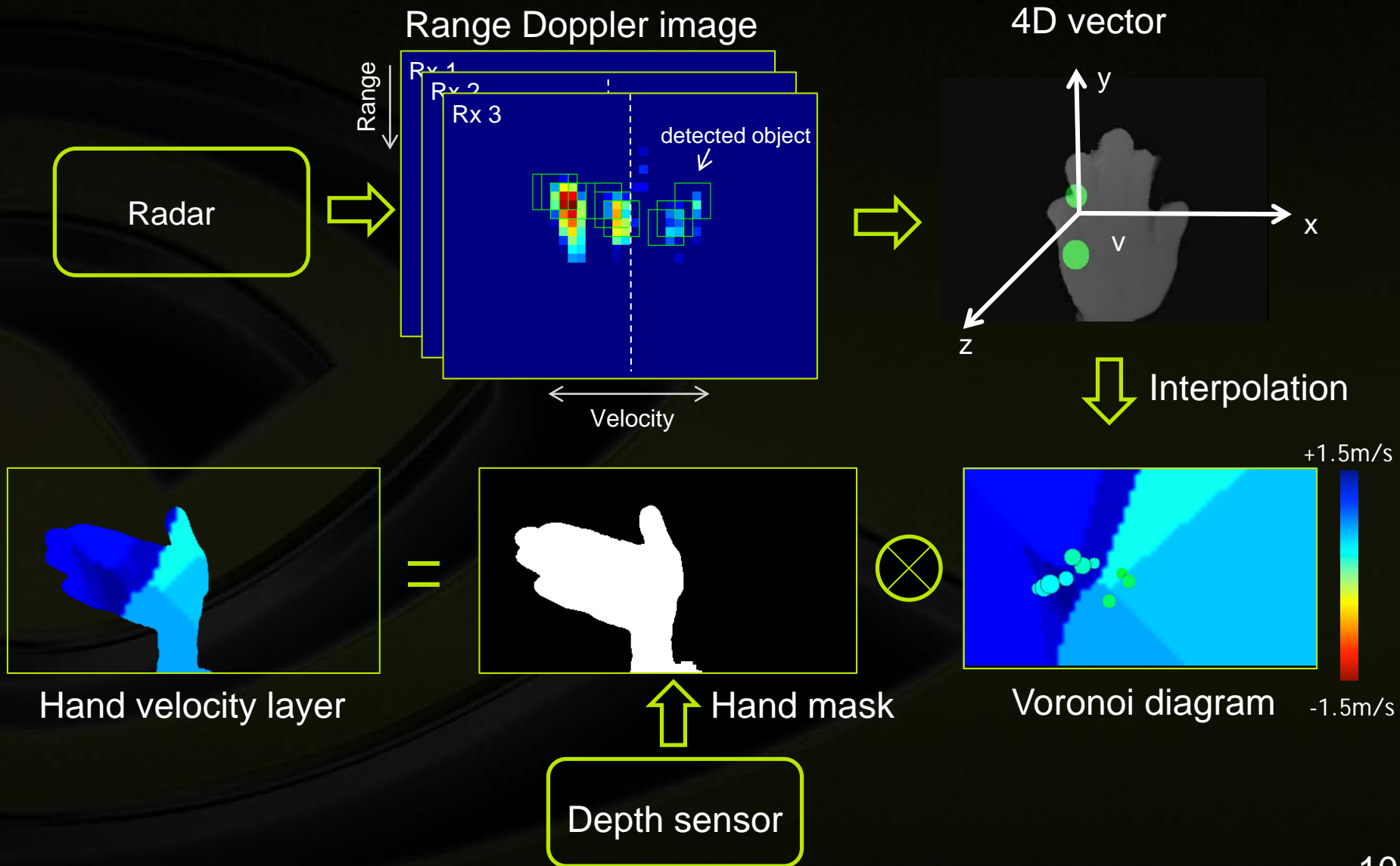
Microcontroller
Analog circuits
Radar chips
Rx2
Rx1
Tx
Rx3
4 antennas

radar prototype

Rx – receiver
Tx – transmitter

Rx1
Rx2
Rx3
Tx1
Azimuth
Elevation

9

# Radar pipeline

Range Doppler image

Range

Rx 1

Radar →

Velocity

# Radar pipeline



Range Doppler image

# Radar pipeline

Range Doppler image

Range

Radar →

Rx 1

Rx 2

Rx 3

detected object

Velocity

# Radar pipeline

Range Doppler image

4D vector

Radar

Range

Rx 1
Rx 2
Rx 3

detected object

Velocity

y

x

z

v

Interpolation

+1.5m/s

=

⊗

Hand velocity layer

Hand mask
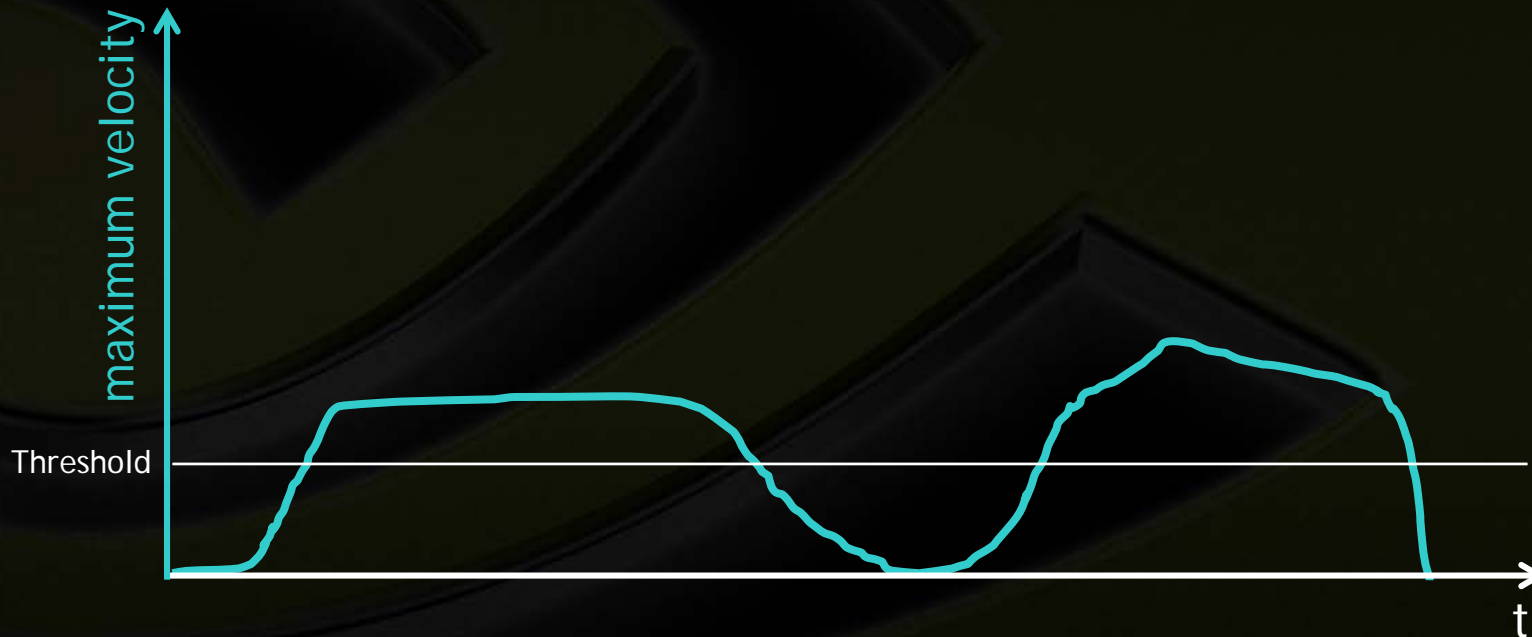
Voronoi diagram

-1.5m/s

Depth sensor

10

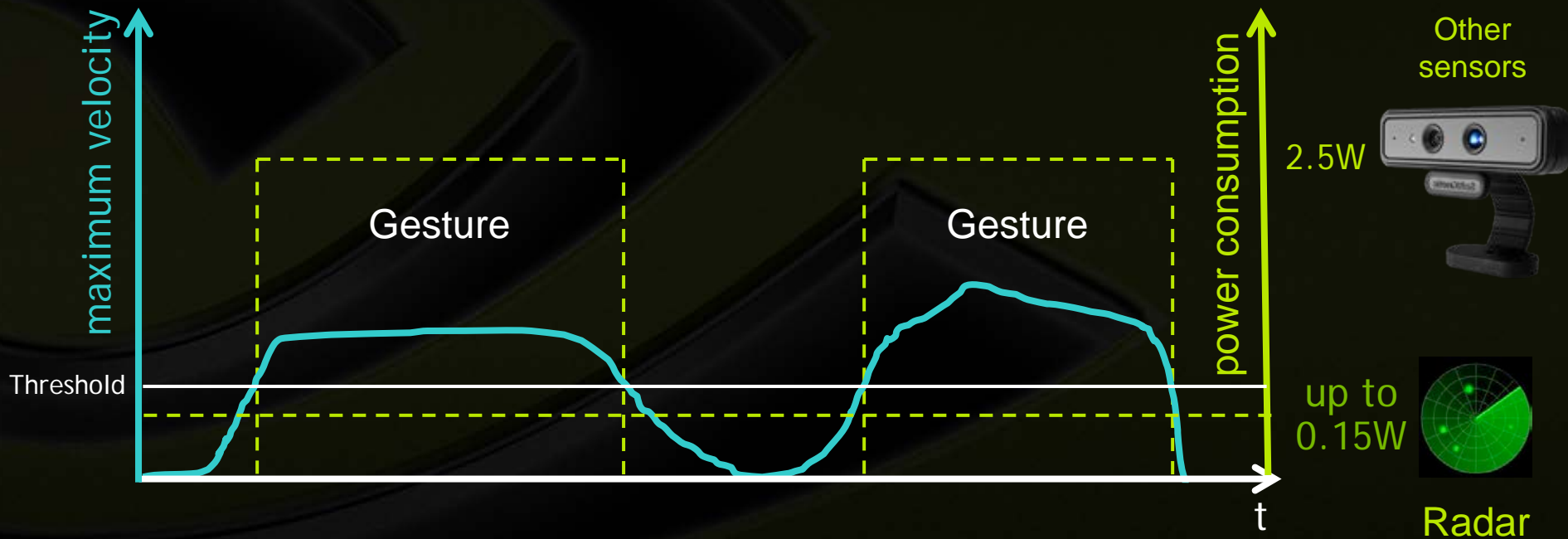# Radar sensor

# Segmentation

- Performed by radar

- Gesture detected when _maximum velocity > threshold_

- Assumptions:

  - Hand stationary between gestures

  - Gesture duration 0.3 - 3s



maximum velocity

Threshold

t

# Power efficiency

- Concept: High power sensors can be switched ON only during gesture detected by radar

- 16x power efficiency

# Classification

- Each frame consists of 3 channels: Intensity, Depth and Velocity

- Each channel is down sampled 32x32 pixels
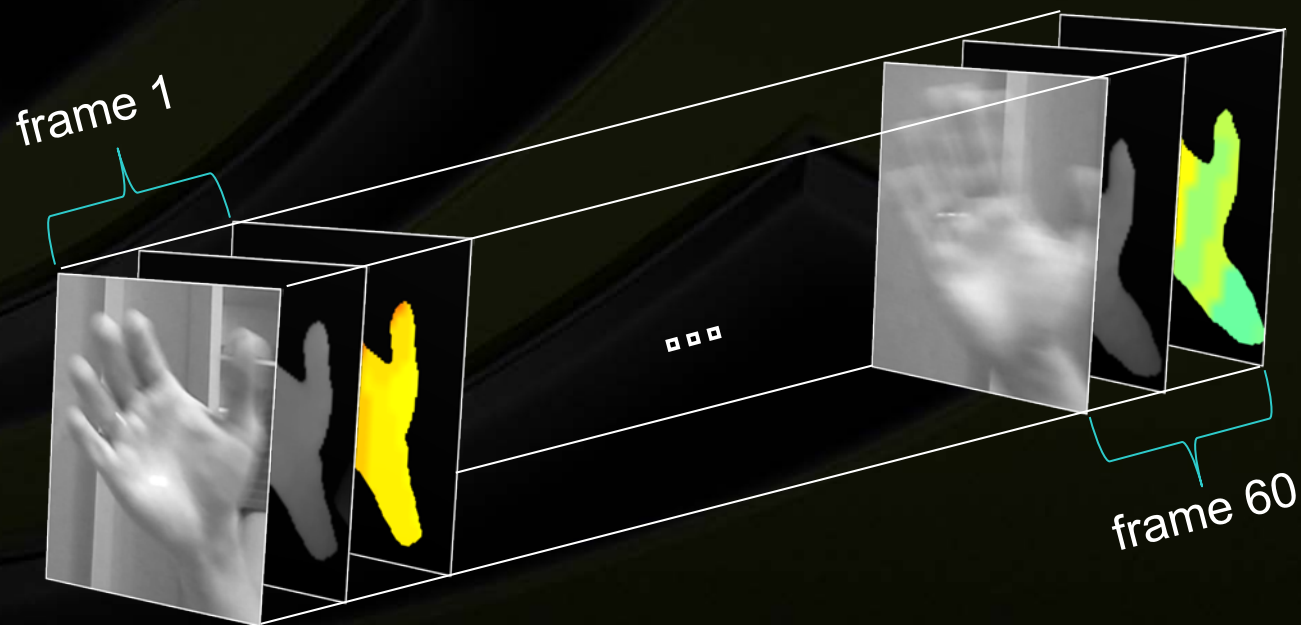


32

32

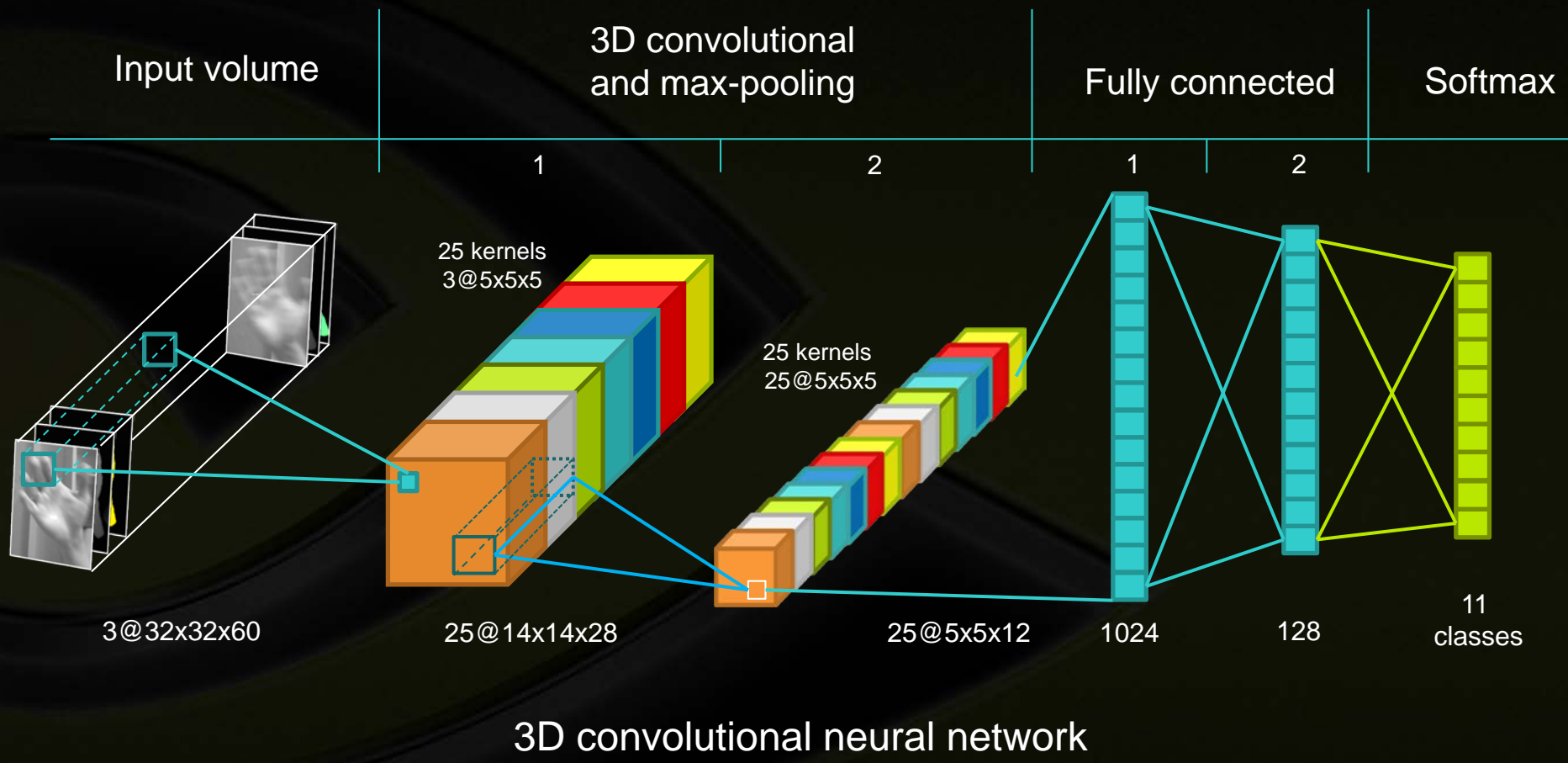Intensity

32

32

Depth

32

32

Velocity

# Classification

- Each frame consists of 3 channels: Intensity, Depth and Velocity

- Each channel is down sampled 32x32 pixels

- Segmented gesture is interpolated to 60 frames

# Classification

- Each frame consists of 3 channels: Intensity, Depth and Velocity

- Each channel is down sampled 32x32 pixels

- Segmented gesture is interpolated to 60 frames



frame 1

frame 60

# Gesture classifier



Input volume

3D convolutional
and max-pooling

Fully connected

Softmax

1

2

1

2

25 kernels
3@5x5x5

25 kernels
25@5x5x5

3@32x32x60

25@14x14x28

25@5x5x12

1024

128

11
classes

3D convolutional neural network

14

# Data collection

- 10 experiments (sessions)
- 10 gesture types + random gesture
- 10-20 repetitions
- 3 persons

- Set-ups:
  - Indoor simulator
  - Outdoor, parked car

- Lighting conditions:
  - Day, shadow
  - Day, Sun
  - Evening
  - Night

- Total 1714 gestures

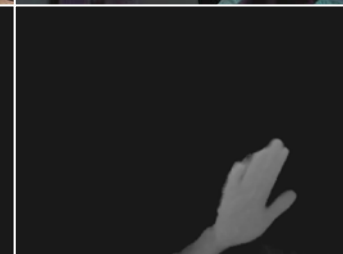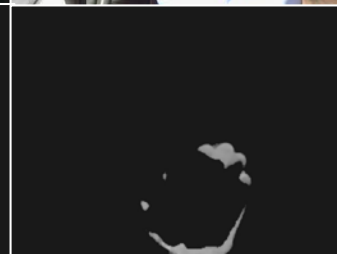Indoor car simulator



Outdoor car

# Data collection

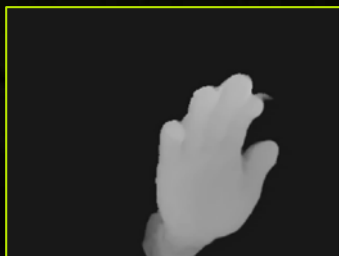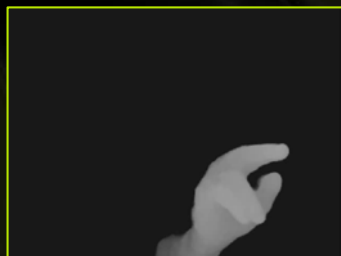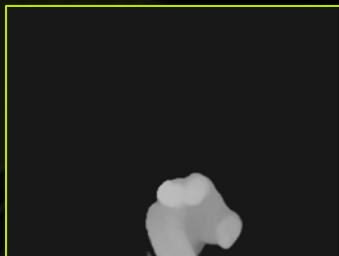# Gesture classes



PALM: left, right, up, down

SWIPE: left, right

SHAKE

ROTATION: CW, CCW

CALL

# Results
## Leave one session out



Classification error, %

- 39,90% — O
- 9,10% — D
- 10,90% — R
- D+O
- R+D
- R+O
- R+D+O

D – depth        O – optical        R - radar

# Results
## Leave one session out

# Results
**Leave one session out**



Classification error, %

- 39,90% — O
- 9,10% — D
- 10,90% — R
- 7,90% — D+O
- 8,30% — R+D
- 7,40% — R+O
- 5,90% — R+D+O

D – depth     O – optical     R - radar

# Results
## Leave one session out: different lighting conditions



Classification error, %

6,70% — Night
3,00% — Evening
9,70% — Day (shadow)
20,90% — Day (sunlight)

■ D+R (CNN)

D – depth    O – optical    R - radar

# Results
## Leave one session out: different lighting conditions



Classification error, %

- 🟩 D+R (CNN)
- 🟧 D+R+O (CNN)

6.70%  6,70%  —  Night

3,00%  1,50%  —  Evening

9,70%  8,30%  —  Day (shadow)

20,90%  7,50%  —  Day (sunlight)

D – depth     O – optical     R - radar

# Results

## Leave one session out: comparison with HOG features



Classification error, %

22,20%

6,70%

1,50% 2,45%

13,00%

20,90%

8,30%

7,50%

Night

Evening

Day (shadow)

Day (sunlight)

D+R+O (CNN)

D+O (HOG*)

D – depth      O – optical      R - radar

*Ohn-Bar and Trivedi, *IEEE Trans. on Intelligent Transportation Systems*, 2014.

# Results

Classification error for leave one subject out classification

| D+R+O (CNN) | D+O (HOG*) |
|:---:|:---:|
| **24.90%** | 48.2% |

D – depth          O – optical          R - radar

Solutions:

- Use more subjects to train

- Perform biometric registration of the system for a new user

*Ohn-Bar and Trivedi, *IEEE Trans. on Intelligent Transportation Systems*, 2014.

# Demo
## Gesture classification in parked car

# Conclusions

- A multi-sensor system improves accuracy and robustness to lighting

- CNN combines sensors effectively

- Proposed approach outperforms feature-based SOTA

- Using radar lowers power consumption, allows efficient gesture segmentation, and improves classification accuracy

THANK YOU
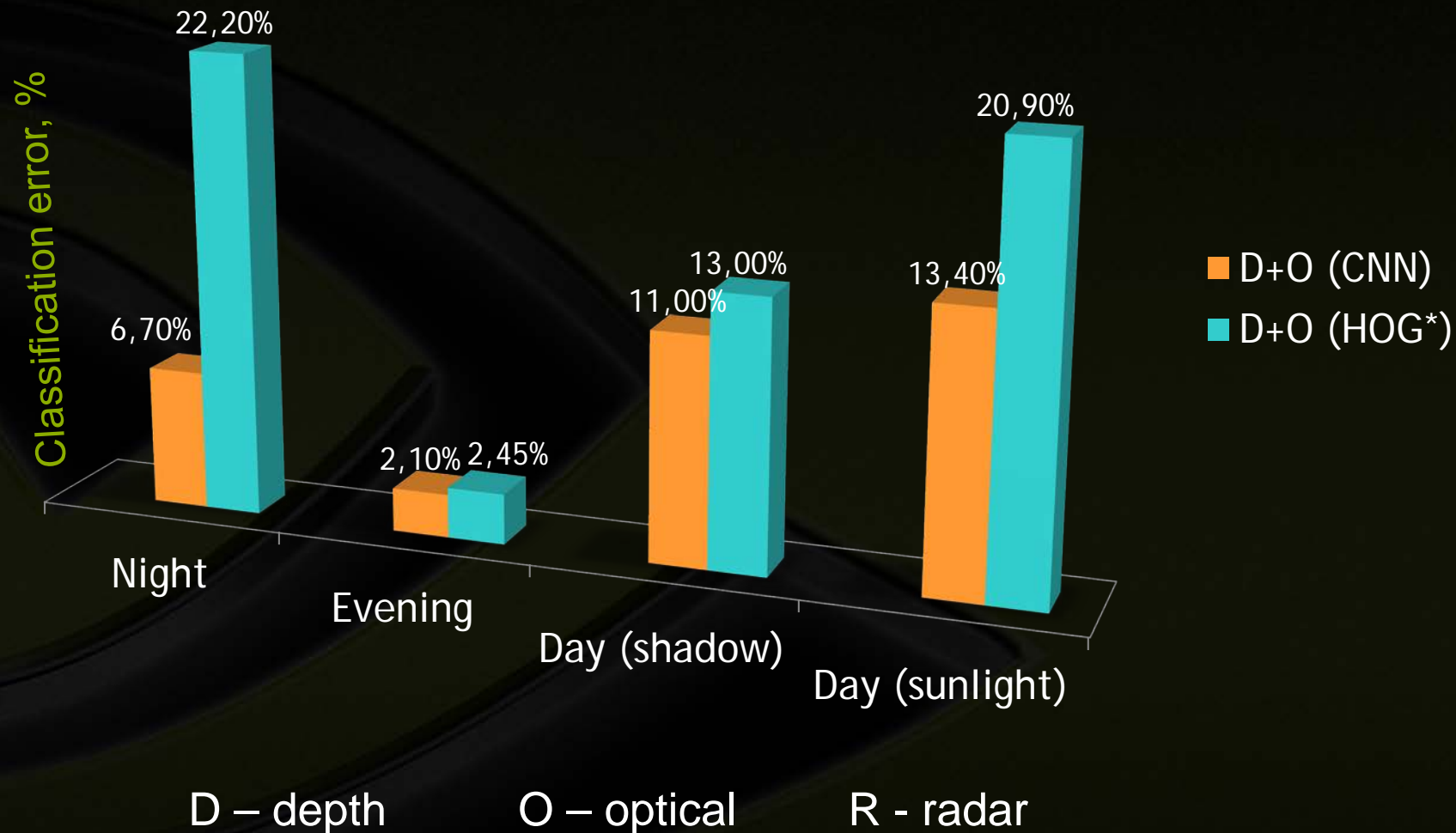QUESTIONS?

# Demo

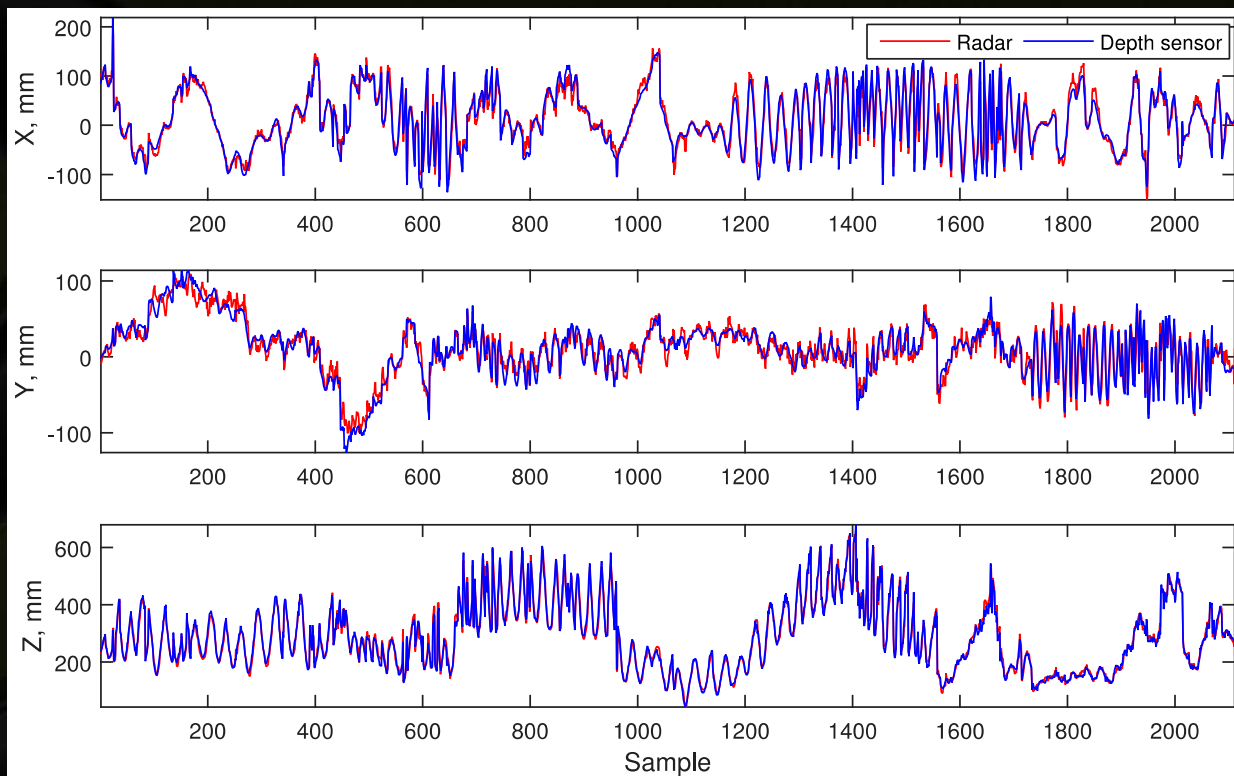**Gesture classification in simulator**



Detected

# Results
## Leave one session out: comparison with HOG features

Classification error, %

22,20%

20,90%

13,00%
11,00%

13,40%

6,70%

2,10% 2,45%

■ D+O (CNN)
■ D+O (HOG*)

Night

Evening

Day (shadow)

Day (sunlight)

D – depth          O – optical          R - radar

*Ohn-Bar and Trivedi, *IEEE Trans. on Intelligent Transportation Systems*, 2014.

# Calibration

Depth and Radar sensors are calibrated
assuming a linear transformation model
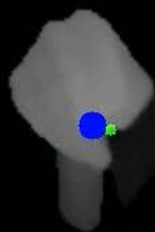
# Radar pipeline

Depth image

Velocity image

+1.5m/s

Detected

-1.5m/s

Green circles – points detected by radar
Blue point – hand position estimated from the radar