# Deep Learning for Object Detection in Video Surveillance

Michelangelo Fiore & Florian Matusek

**KiwiSecurity**
AUTOMATING VIDEO SURVEILLANCE

Number of video surveillance cameras world-wide

300m+

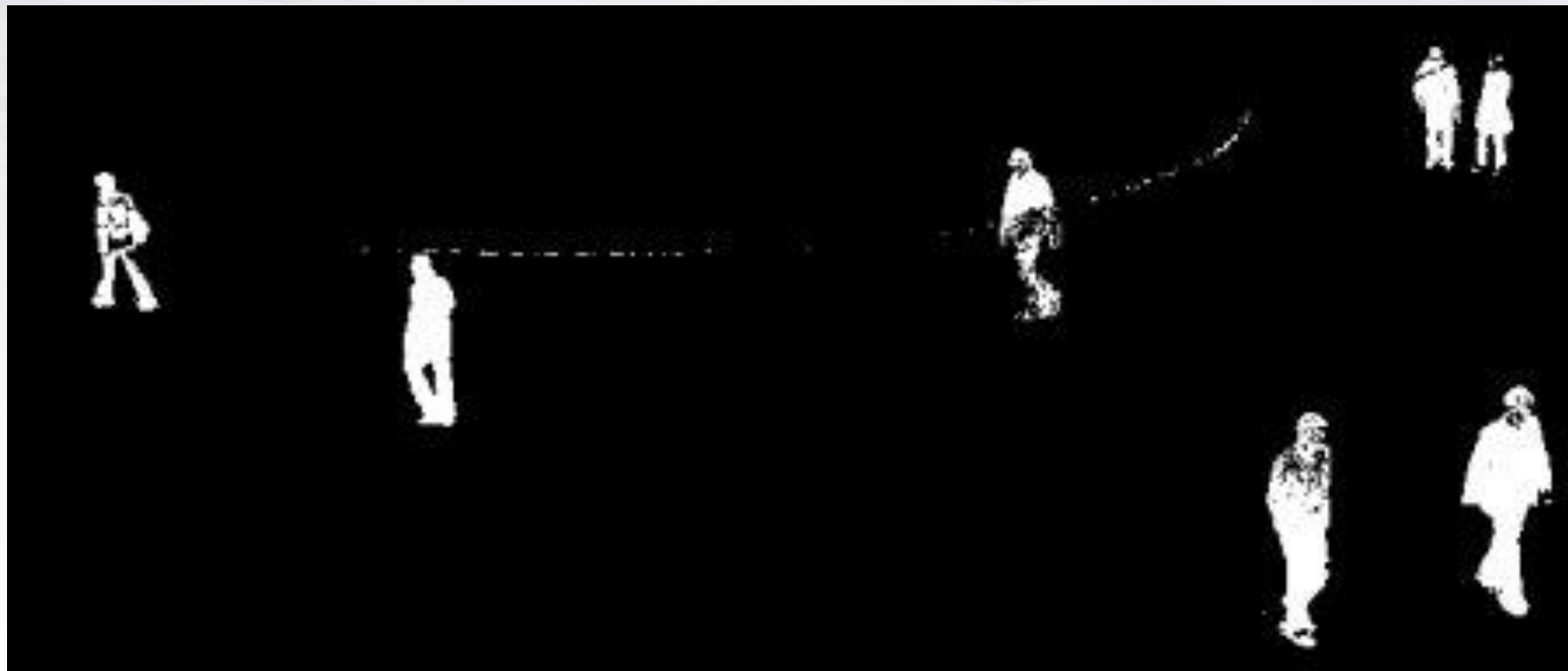**Incidents missed by security operators after 20 minutes**

95%

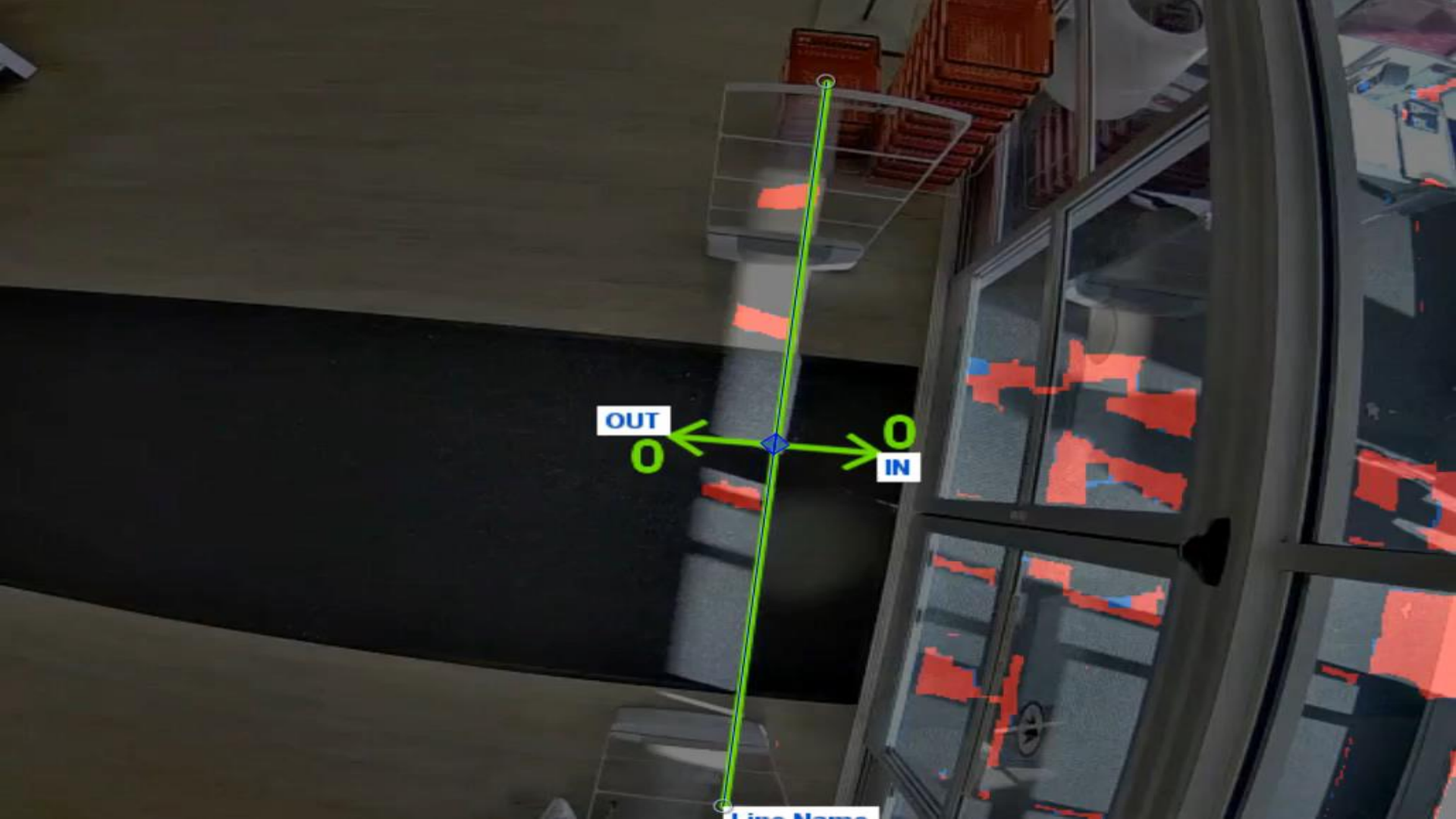# So where do we go from here?

# Background subtraction

Kiw iSecurity – Automating Video Surveillance

# Background subtraction

Kiw iSecurity – Automating Video Surveillance

Some Pitfalls

OUT
0
0
IN

Line Name
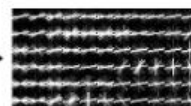
# Many Parameters

# Hard to Tune

**Too Much**

# Choosing the Right Features
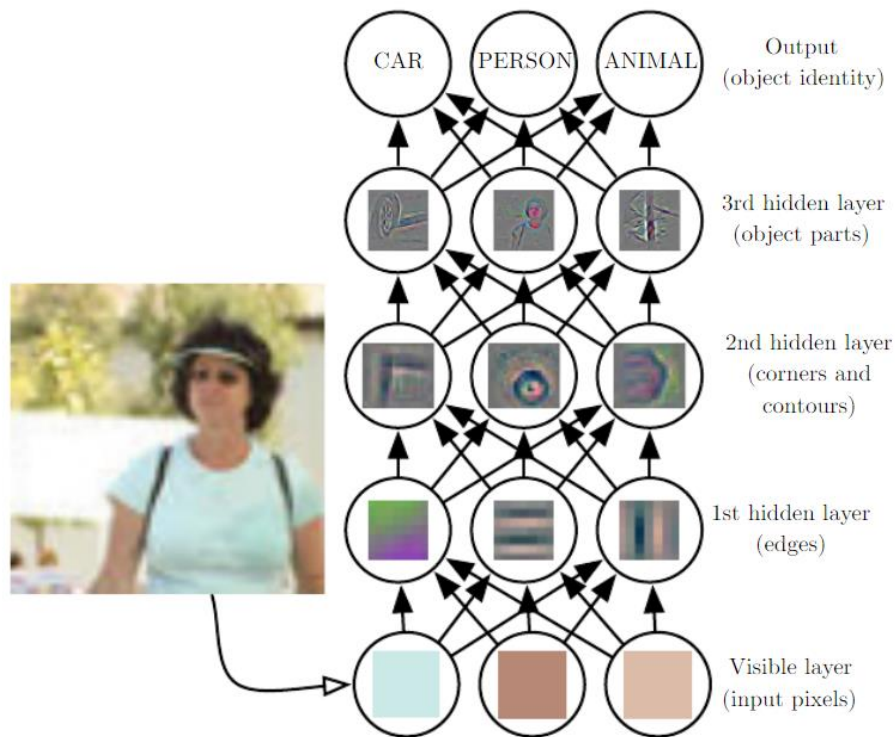

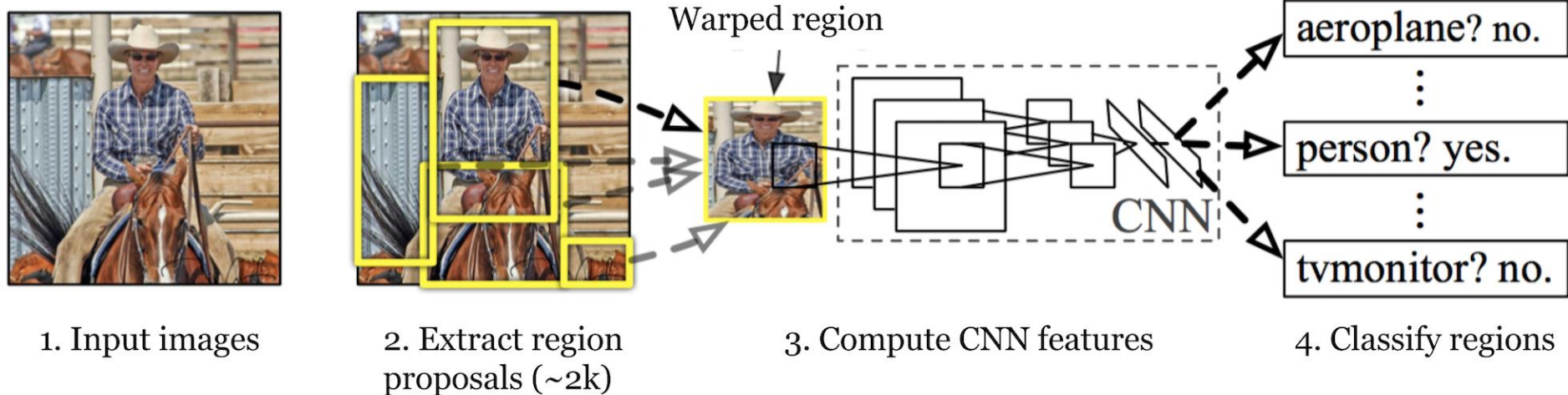
Car



Car Detection → HOG Features → Our Visualization

C. Vondrick, A. Khosla, T. Malisiewicz, A. Torralba. "HOGgles: Visualizing Object Detection Features" *International Conference on Computer Vision* (ICCV), Sydney, Australia, December 2013.

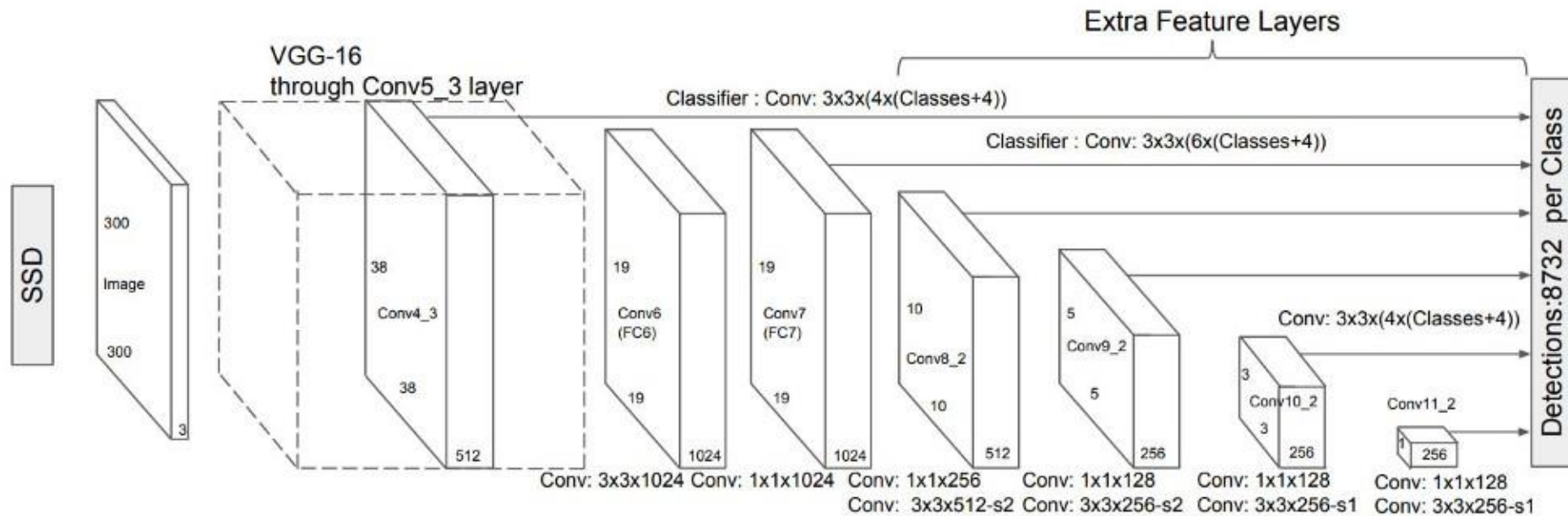# Feature Extraction With Deep Learning

# Object Detection with Deep Learning: Region Based Approach



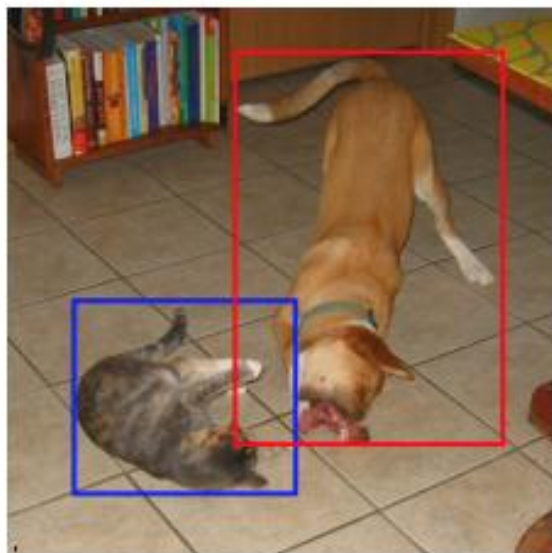1. Input images    2. Extract region proposals (~2k)    3. Compute CNN features    4. Classify regions

Girshick, Ross, et al. "Rich feature hierarchies for accurate object detection and semantic segmentation." *Proceedings of the IEEE conference on computer vision and pattern recognition.* 2014.

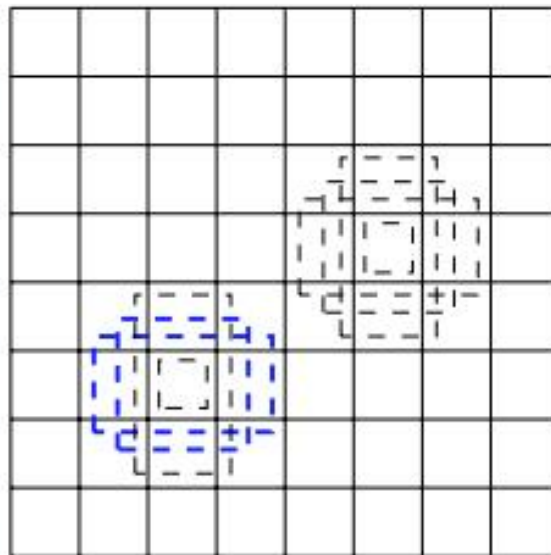# Object Detection with Deep Learning: Single Shot Approach

Liu, Wei, et al. "Ssd: Single shot multibox detector." *European conference on computer vision*. Springer, Cham, 2016
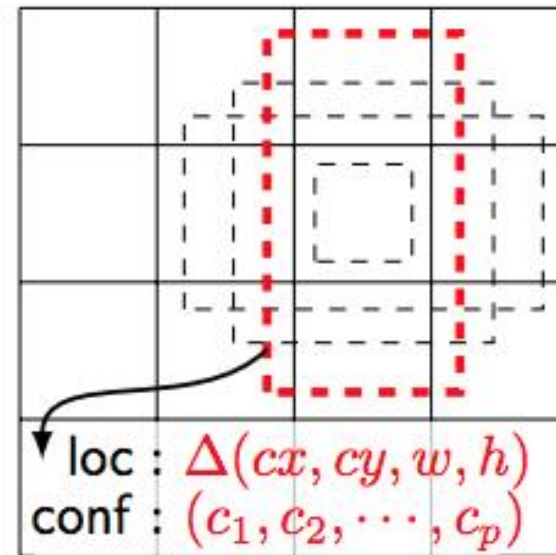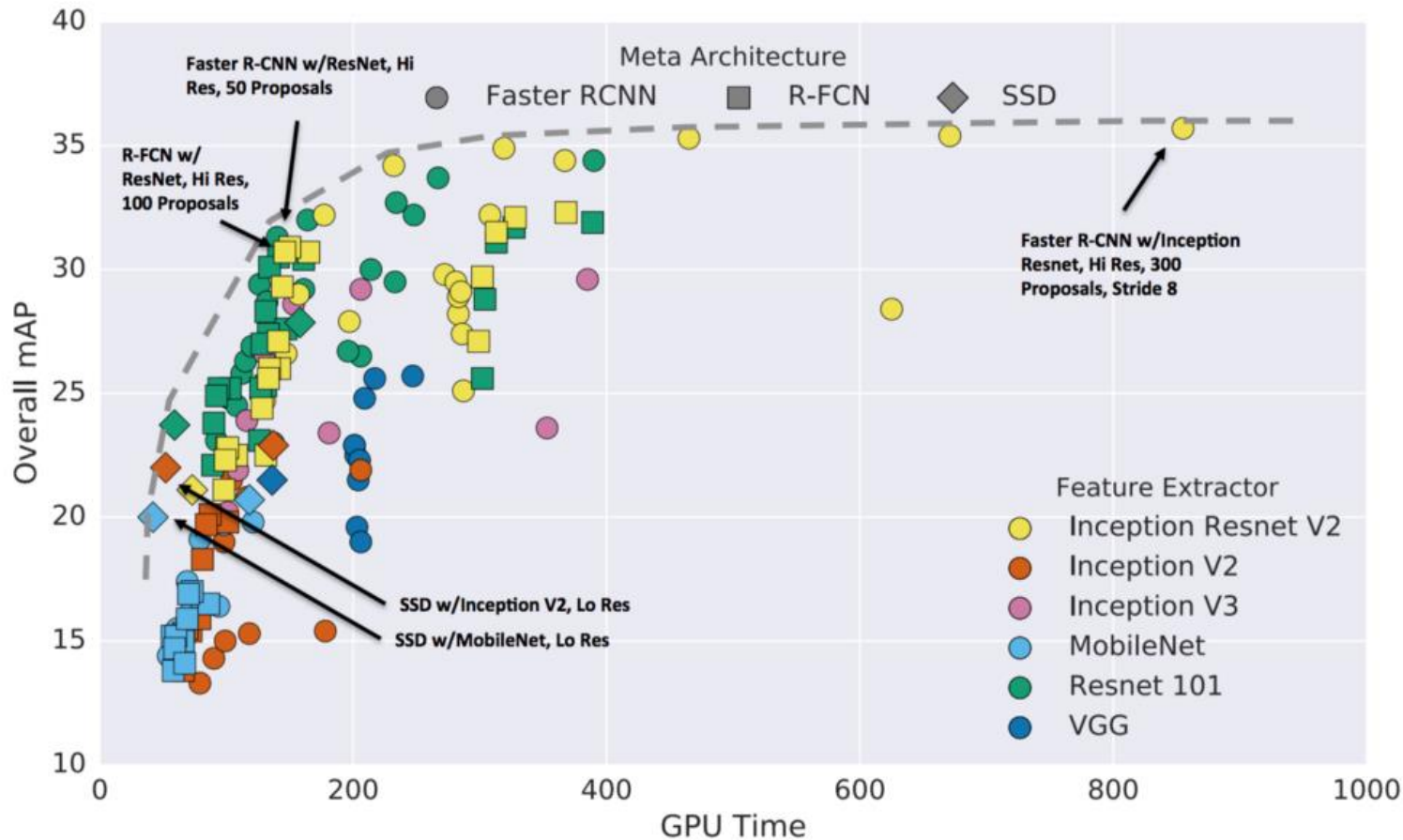
# Detection at Different Scales



(a) Image with GT boxes  (b) 8 × 8 feature map  (c) 4 × 4 feature map
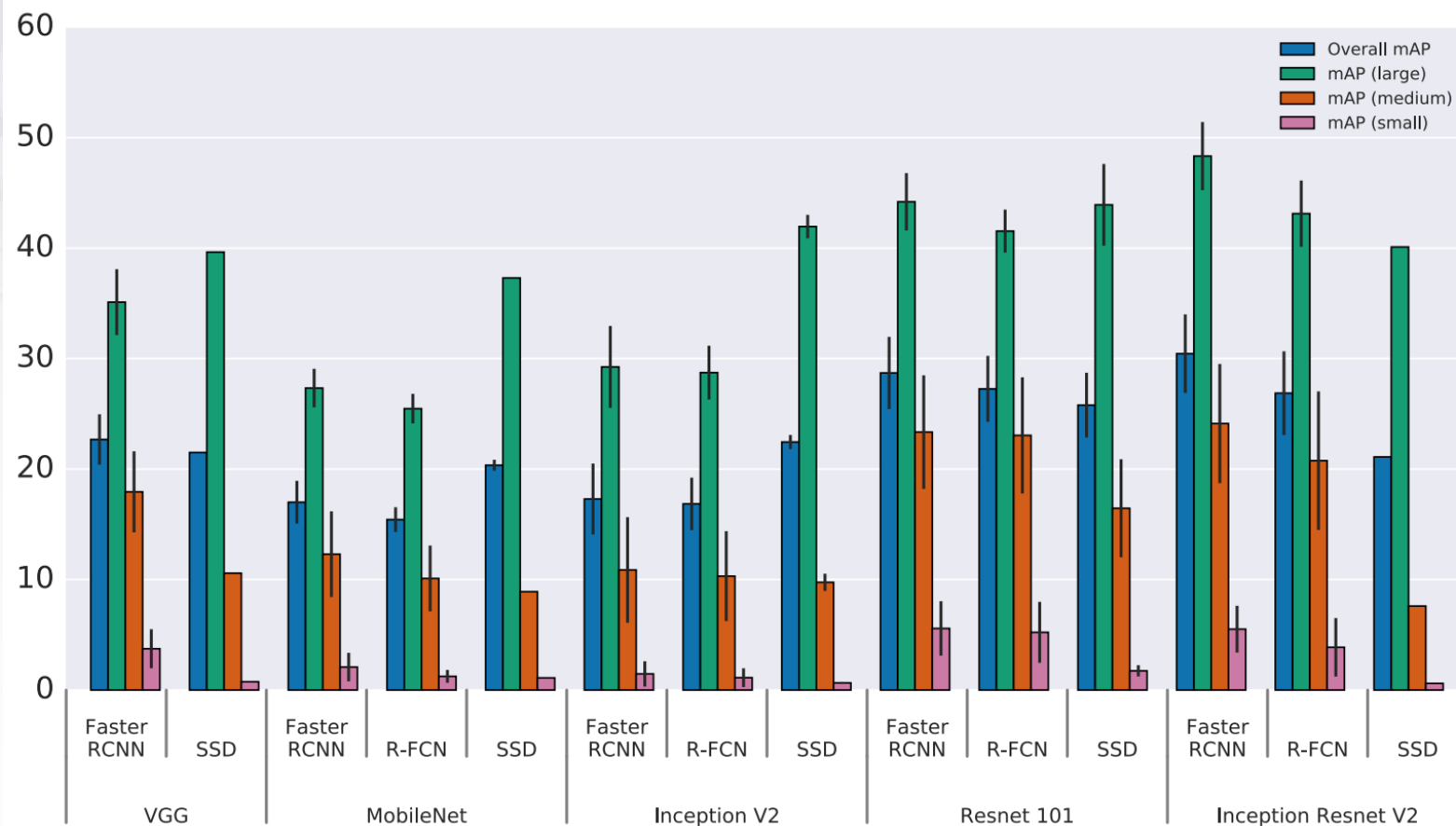
$$\text{loc} : \Delta(cx, cy, w, h)$$
$$\text{conf} : (c_1, c_2, \cdots, c_p)$$

# Choosing a Detector

**Feature Extractors**

**+**

**Meta-Architecture**

| VGG | | SSD |
| Mobilenet | | R-CNN |
| Resnet | | |
| Inception | | |

Huang, Jonathan, et al. "Speed/accuracy trade-offs for modern convolutional object detectors." *IEEE CVPR*. 2017.

# Understanding Your Requirements

Kiw iSecurity – Automating Video Surveillance

# Our Frameworks: Apache MXNet

| OS and Language support |  |
|---|---|
| Documentation | Good API and GLUON |
| Models Support | Many Object Detection Models |
| Performances | Fast Performances on Training and Inference |

https://mxnet.apache.org

mxnet

# Our Frameworks: Intel OpenVINO



https://software.intel.com/en-us/openvino-toolkit

# Our Detector

VGG16

Mobilenet

SSD

Kiw iSecurity – Automating Video Surveillance

# Building a Dataset

Kiw iSecurity – Automating Video Surveillance

# People Detection Dataset

# Some Results

# Improving the Performances
## Reducing the Model Complexity

Kiw iSecurity – Automating Video Surveillance

# Model Pruning



**yes**

**Pruning Percentage**

**Evaluate Filters** → **Remove Chosen Filters** → **Finetune to desired accuracy** → **Continue**

**no**

**Possible Criteria: L2 Norm**

Molchanov, Pavlo, et al. "Pruning convolutional neural networks for resource efficient inference." (2016).

# Pruning Results: VGG16

**Greatly Reduced Model Size**
Some layers cut by 75%
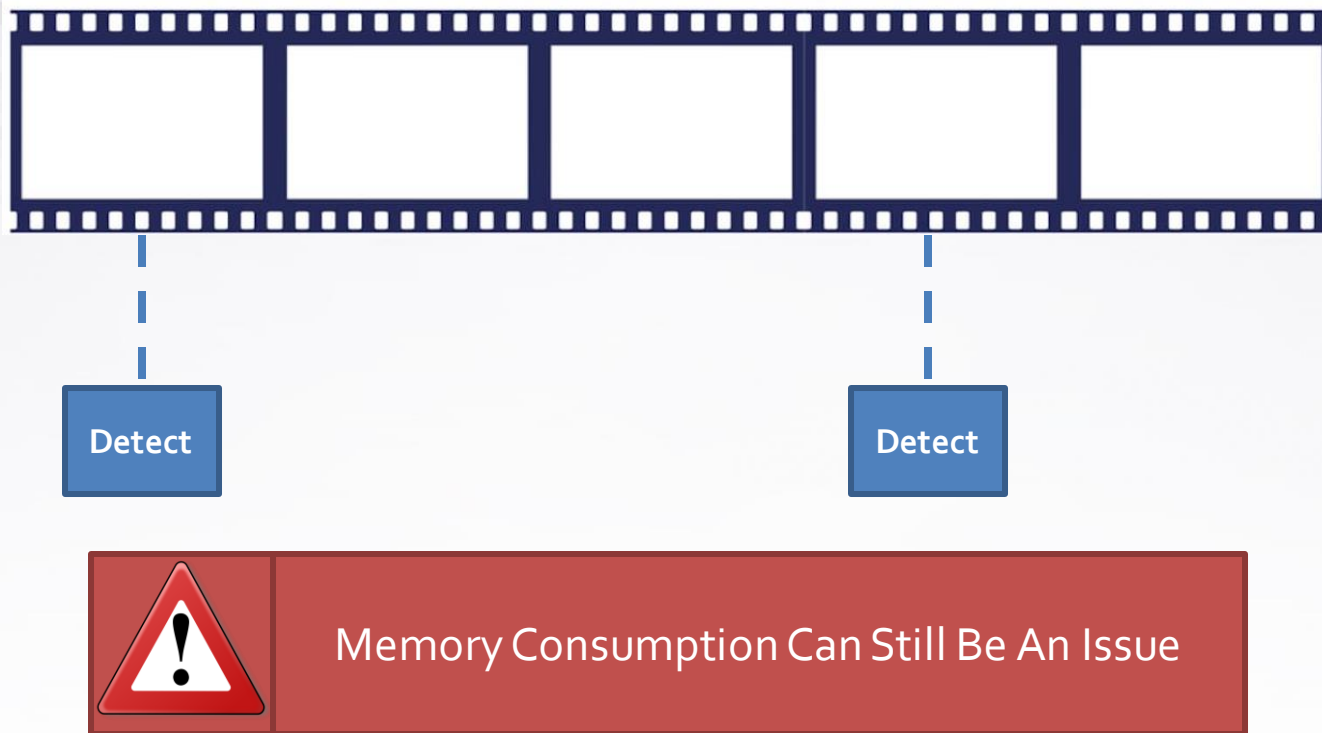
**Accuracy Decreased by 0.005**

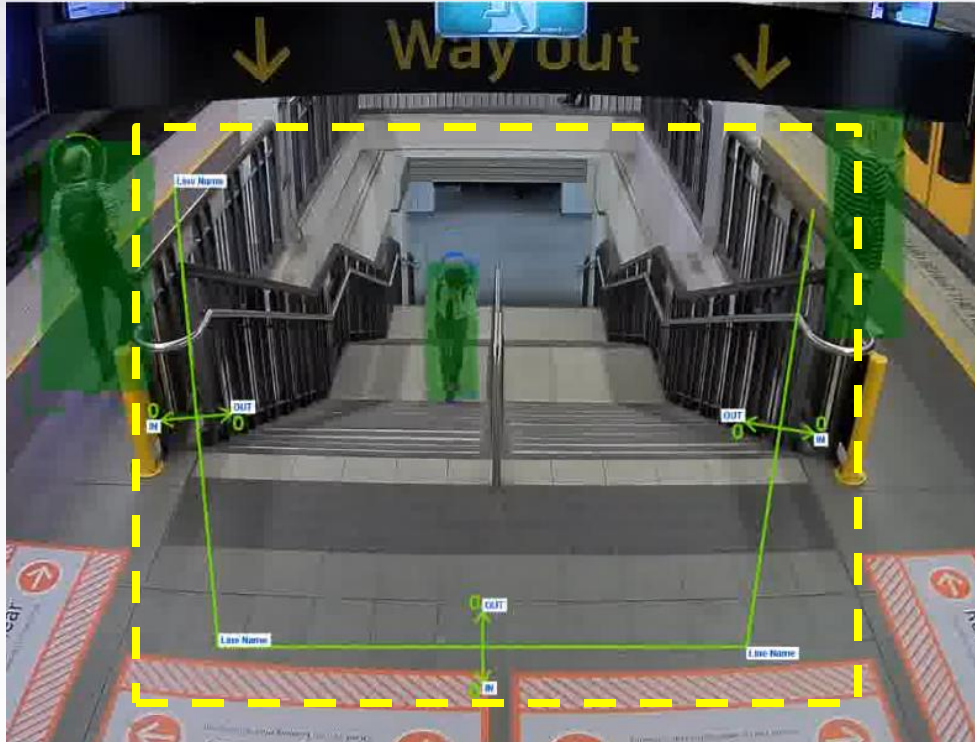**More than Doubled FPS**
30 vs 70

All experiments where done on an Intel Core i7-7800X CPU with a Nvidia Quadro P2000 GPU

# Improving the Performances
## Do We Need to Detect at Every Frame?

Detect

Detect

Memory Consumption Can Still Be An Issue

# Improving the Performances Reducing the Detection Area

| 1 | Reduce Distorsion |
|---|---|
| 2 | Increase Object Size |

Significant Performance Increase (6%)

If  you don't need to beat ImageNet, **don't try to**

Use **every trick** you can

# Future Plans

Explore other domains

Create our own feature extractor

Introduce temporal information in the model

# Interested?

Michelangelo Fiore
m.fiore@kiwisecurity.com

Florian Matusek
f.matusek@kiwisecurity.com


www.kiwisecurity.com


**PS:** We're hiring!