



Università  
della  
Svizzera  
italiana



37th IEEE/ACM International  
Conference on Automated  
Software Engineering

Software Institute 

# ThirdEye: Attention Maps for Safe Autonomous Driving Systems



*Andrea Stocco, USI*



*Paulo J. Nunes, UFPE*



*Marcelo d'Amorim, UFPE*



*Paolo Tonella, USI*

 @tsigalko18

Lane keeping & detection

Object recognition

vehicles

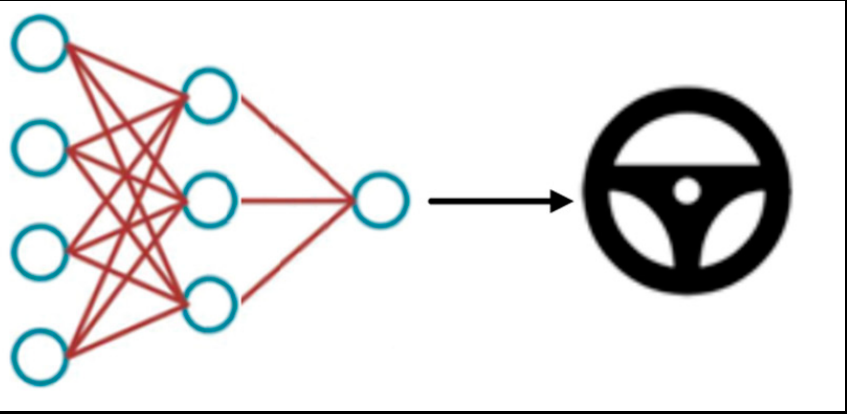
pedestrians

/Autonomous  
/Sensing  
/Communication  
/Battery  
/Navigation  
/Mirrorless  
/Ecology

Self-Driving

48  
mph

**Model Construction  
(training)**



**Data Collection  
(in-field driving)**



**In-field testing  
(evaluation)**



**In-house testing  
(simulation)**



## Steering Angle predicted by the DNN

-0.16933852434158325  
-0.2953818142414093  
-0.2953818142414093  
-0.2953818142414093  
-0.24482464790344238  
-0.24482464790344238  
-0.24482464790344238  
-0.2340604066848755  
-0.2340604066848755  
-0.2340604066848755  
-0.2876757085323334  
-0.2876757085323334  
-0.2876757085323334  
-0.28597092628479004  
-0.28597092628479004  
-0.28597092628479004  
-0.280177503824234  
-0.280177503824234  
-0.280177503824234  
-0.1850987821817398  
-0.1850987821817398  
-0.1850987821817398  
-0.2626234292984009  
-0.2626234292984009  
-0.2626234292984009  
-0.20239685475826263



**Training set cannot contain  
ALL possible driving conditions!**

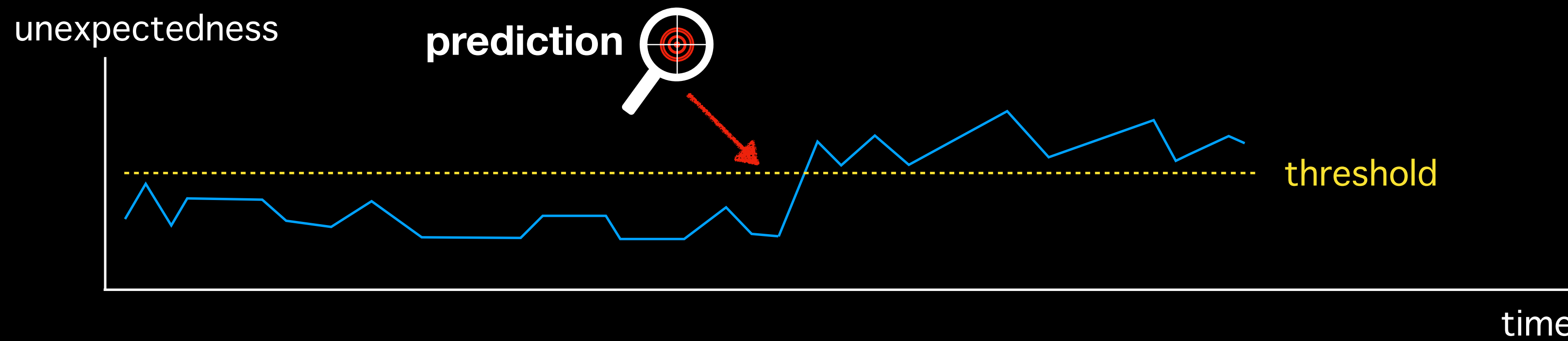
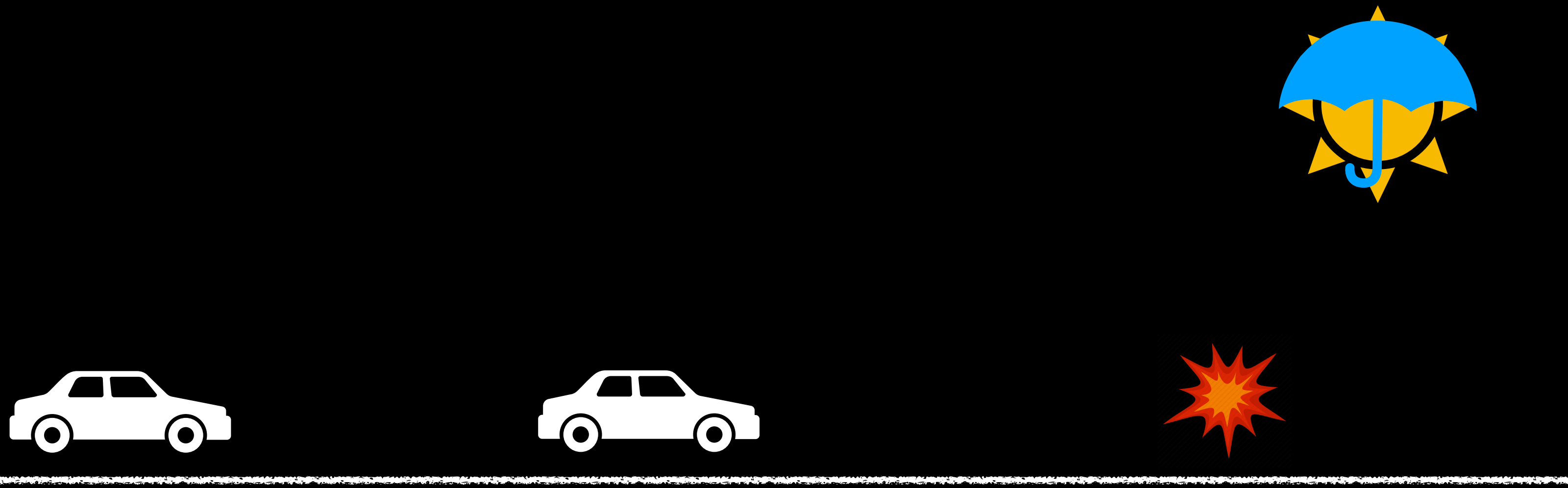
*Can we predict  
unexpected conditions\**



*\* Unexpected condition:*

*Unseen, potentially hazardous and  
misbehaviour-inducing driving condition*

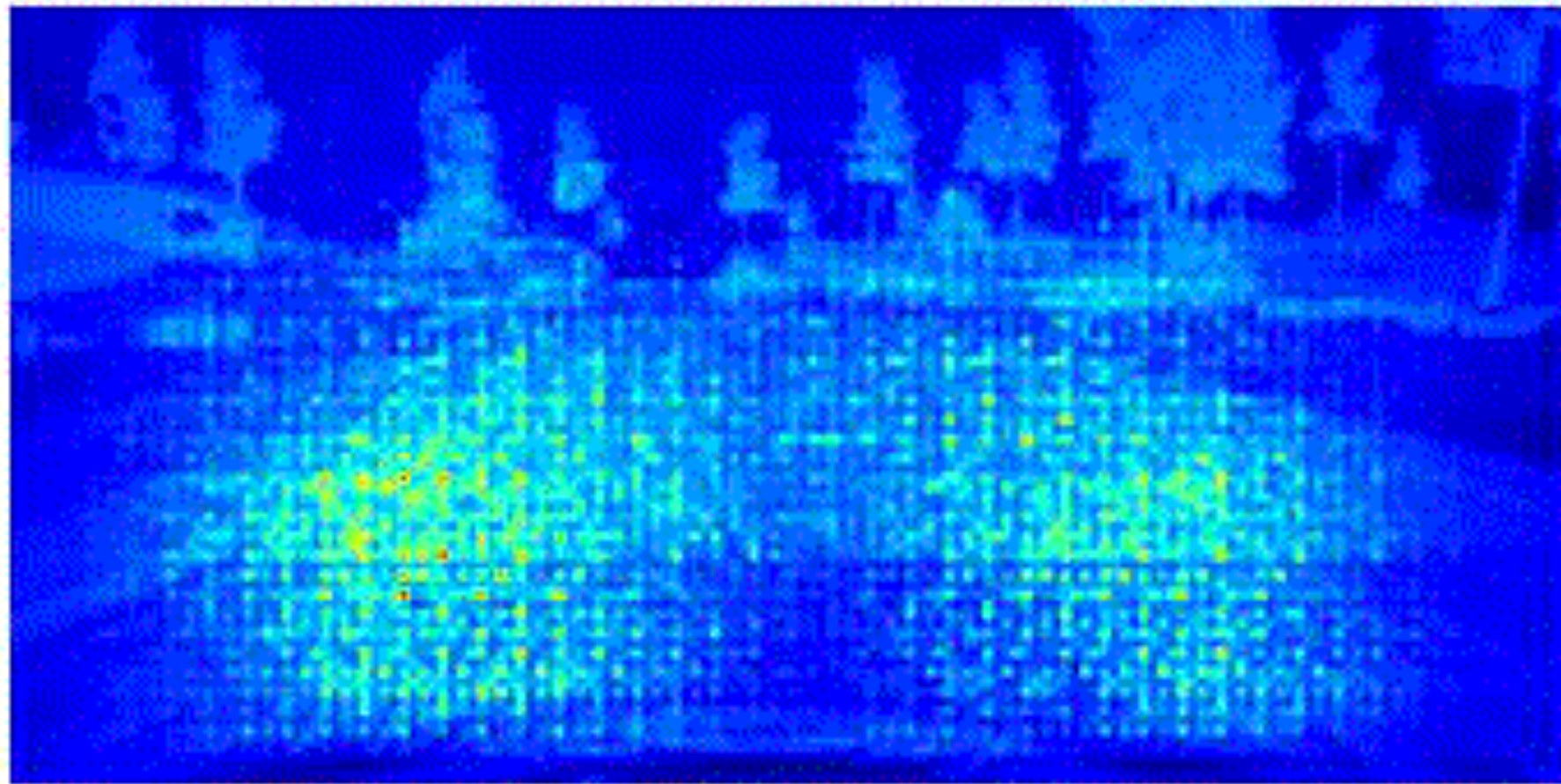
Unexpectedness metric  
Anomaly detection



# XAI for failure prediction

Nominal

XAI Image



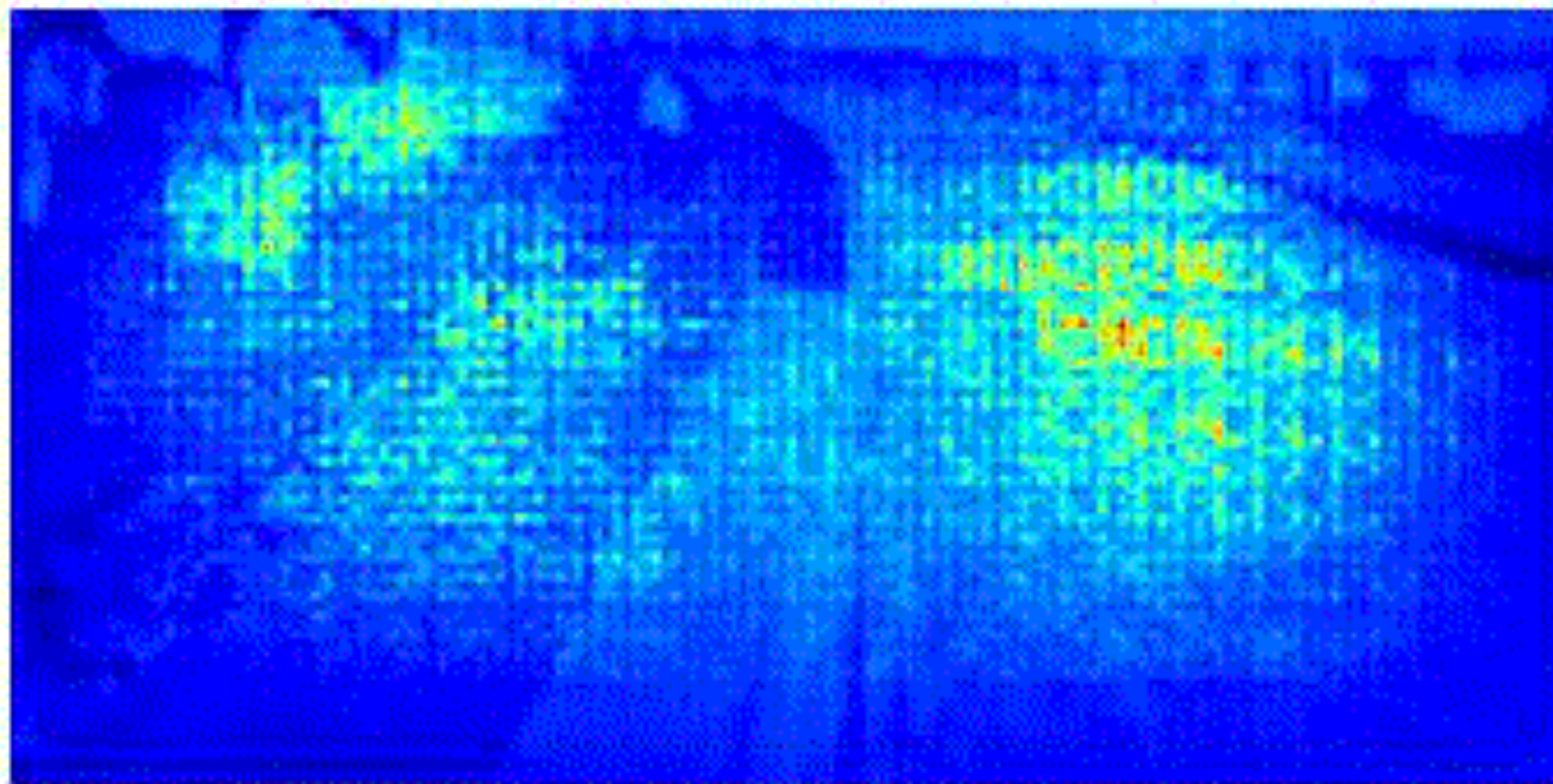
Original Image



# XAI for failure prediction

Uncertain/Unexpected

XAI Image



Original Image



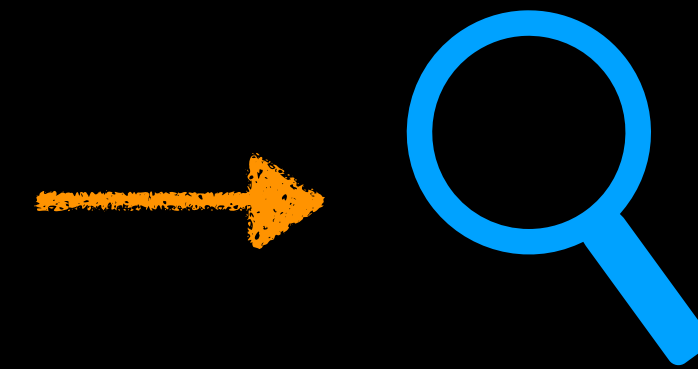




**1. Failure Predictor Training**

**2. Probability distribution fitting**

**3. Treshold estimation**



Trained Predictor



# 1. Failure Predictor Training

## 1.1 Attention Map Generation

## 2. Probability distribution fitting

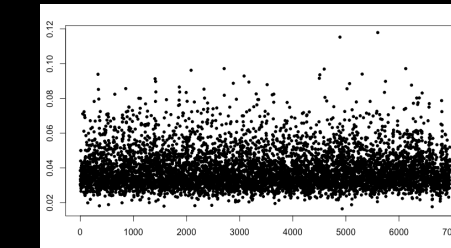
## 3. Treshold estimation

# 1. Failure Predictor Training

1.1 Attention Map Generation

# 2. Probability distribution fitting

2.1 Confidence Score Synthesis



2.2 Gamma distribution Fitting

# 3. Treshold estimation

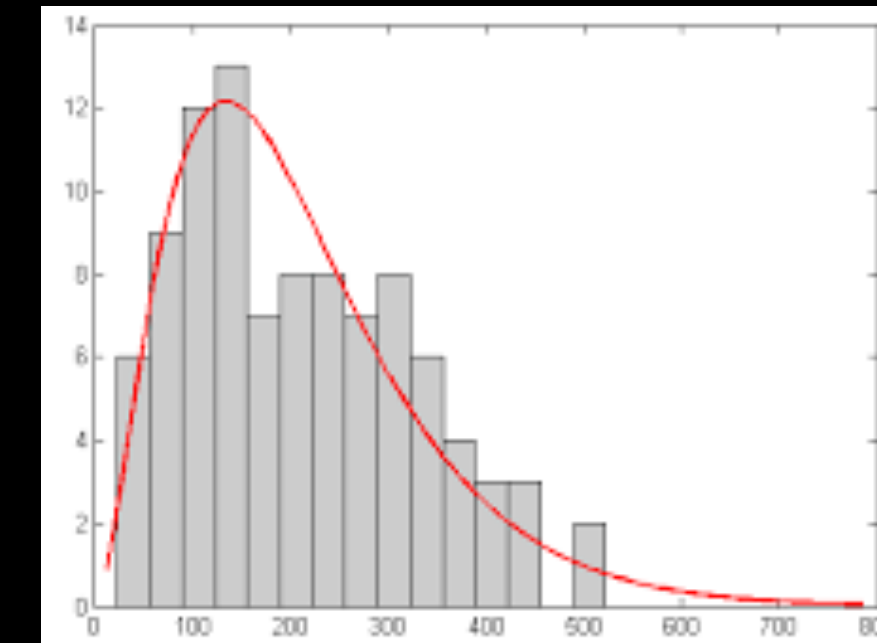
# 1. Failure Predictor Training

1.1 Attention Map Generation

## 2. Probability distribution fitting

2.1 Confidence Score Synthesis

2.2 Gamma distribution Fitting



## 3. Treshold estimation

# 1. Failure Predictor Training

1.1 Attention Map Generation

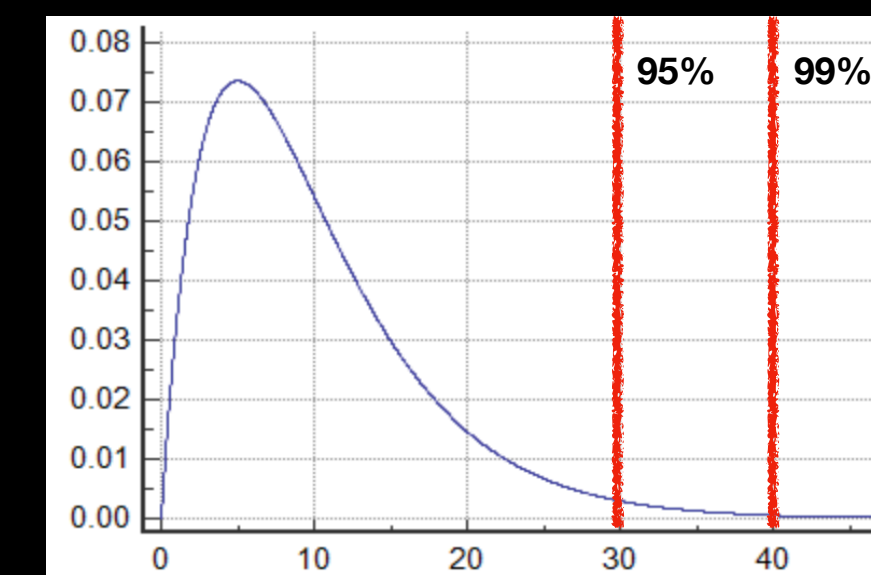
# 2. Probability distribution fitting

2.1 Confidence Score Synthesis

2.2 Gamma distribution Fitting

# 3. Treshold estimation

3.1 w/ Maximum Likelihood Estimation



# 1. Failure Predictor Training

1.1 Attention Map Generation

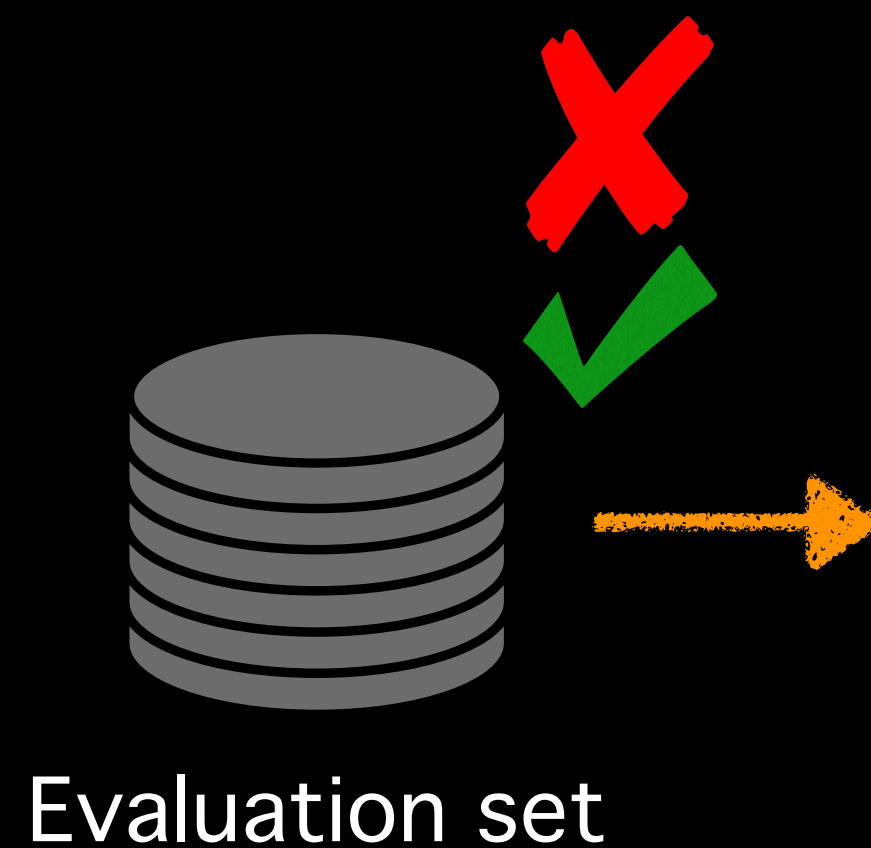
# 2. Probability distribution fitting

2.1 Confidence Score Synthesis

2.2 Gamma distribution Fitting

# 3. Treshold estimation

3.1 w/ Maximum Likelihood Estimation



# 4. Testing

14



# Confidence Score Synthesis

$$\bar{\mathbf{h}} = \frac{1}{WHC} \sum_{i=1, j=1, c=1}^{W, H, C} h_{[i][j][c]}$$

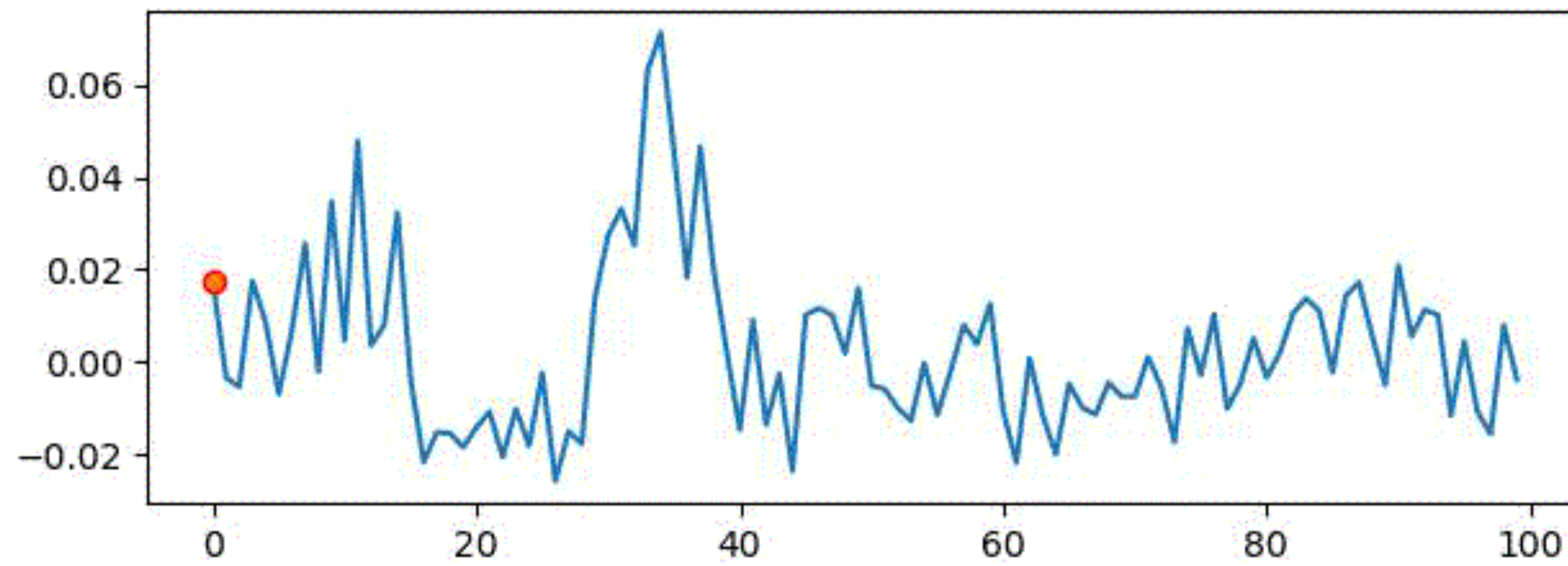
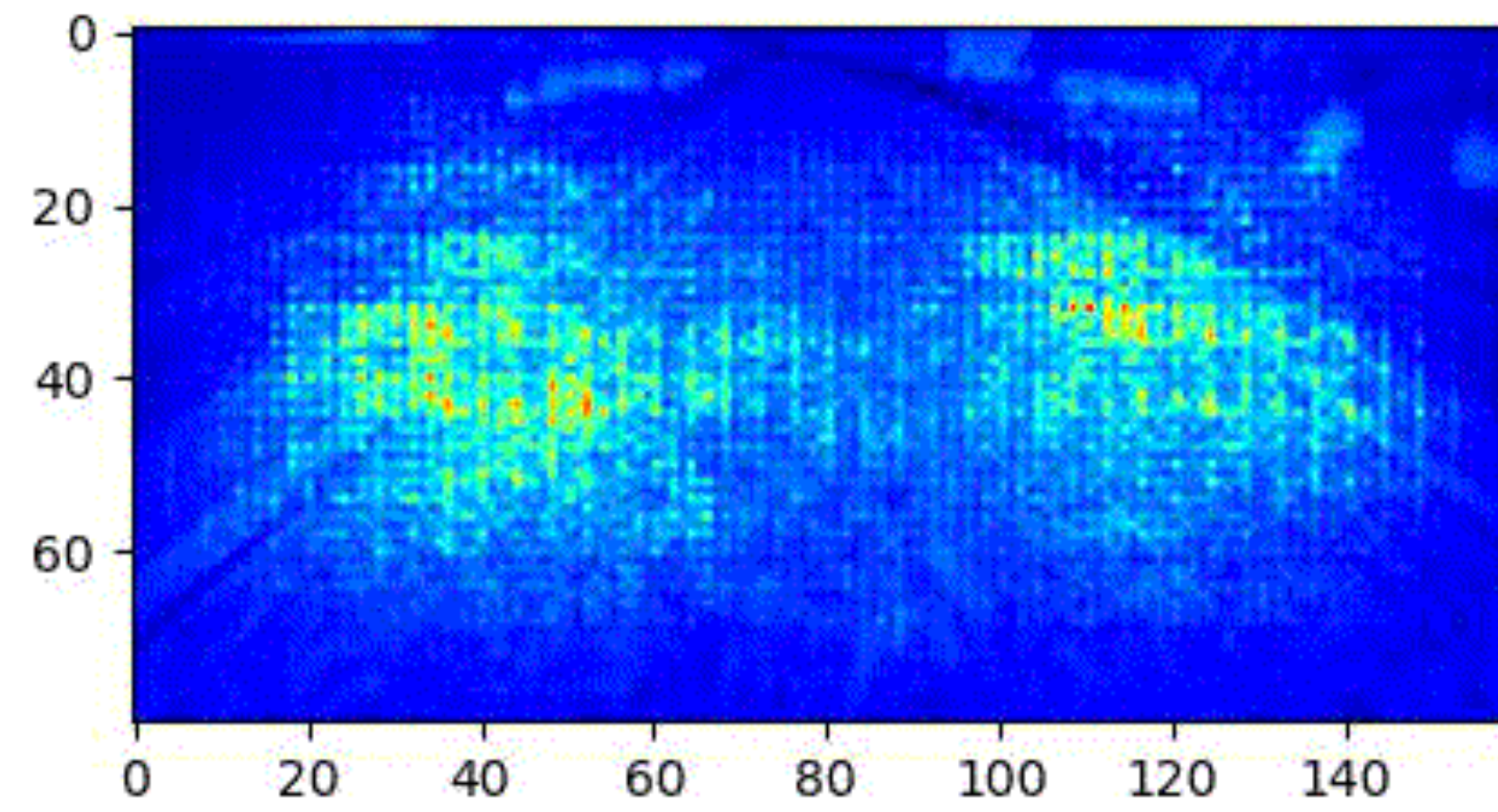
HA = Heatmap  
Average Function

$$\nabla \mathbf{h}_t = \frac{1}{WHC} \sum_{i=1, j=1, c=1}^{W, H, C} h_{t[i][j][c]} - h_{t-1[i][j][c]}$$

HD = Heatmap  
Derivative Function

$$\mathbf{h}_e = \mathcal{L}(\mathbf{h}, \text{dec}(\text{enc}(\mathbf{h})))$$

HRL = Heatmap  
Reconstruction Loss



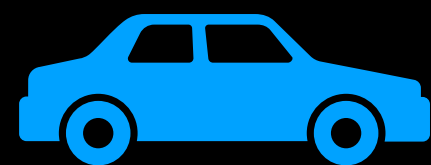


# Experimental Study

Procedure, Main Results



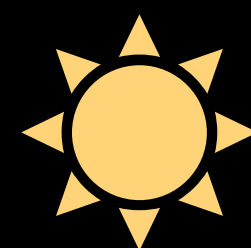
Udacity sim



Dave-2



Predictor (Thirdeye)



## ThirdEye

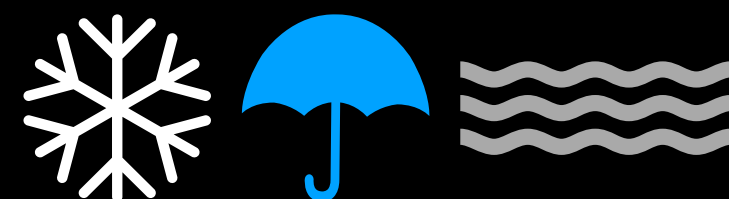
3 confidence scores

- Average
- Derivative
- Rec. Loss

Baseline: SelfOracle (ICSE 2020)

## Two Experiments

External Unknown Scenarios



Internal Uncertain Scenarios



## Metrics

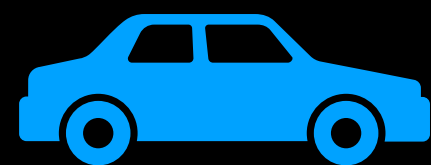
Precision

Recall

F-3



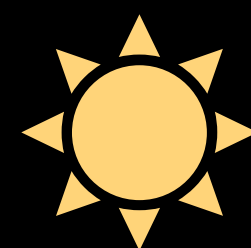
Udacity sim



Dave-2



Predictor (Thirdeye)



## ThirdEye

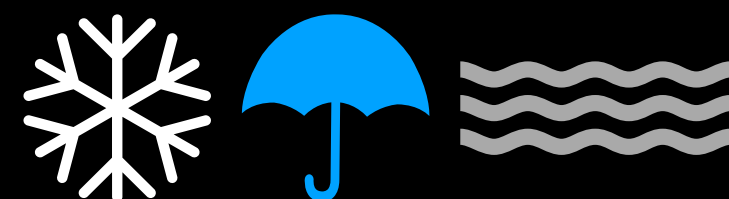
3 confidence scores

- Average
- Derivative
- Rec. Loss

Baseline: SelfOracle (ICSE 2020)

## Two Experiments

External Unknown Scenarios



Internal Uncertain Scenarios



## Metrics

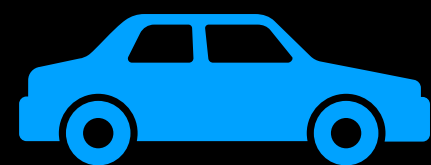
Precision

Recall

F-3



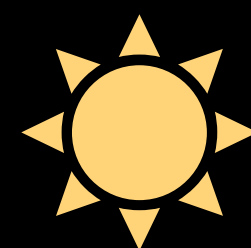
Udacity sim



Dave-2



Predictor (Thirdeye)



## ThirdEye

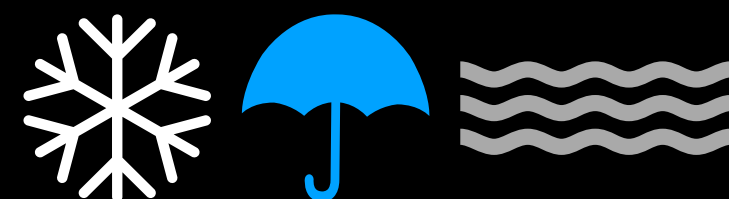
3 confidence scores

- Average
- Derivative
- Rec. Loss

Baseline: SelfOracle (ICSE 2020)

## Two Experiments

External Unknown Scenarios



Internal Uncertain Scenarios



## Metrics

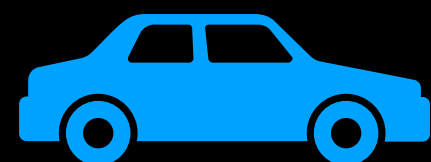
Precision

Recall

F-3



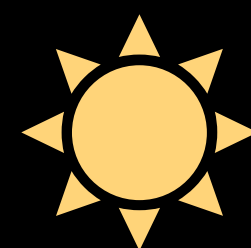
Udacity sim



Dave-2



Predictor (Thirdeye)



## ThirdEye

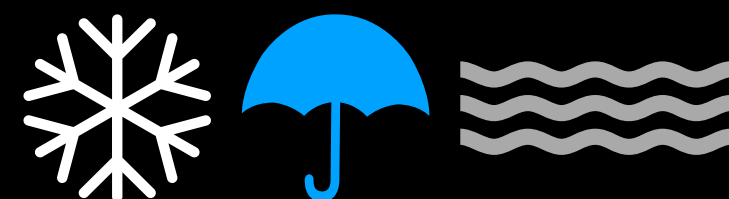
3 confidence scores

- Average
- Derivative
- Rec. Loss

Baseline: SelfOracle (ICSE 2020)

## Two Experiments

External Unknown Scenarios



Internal Uncertain Scenarios

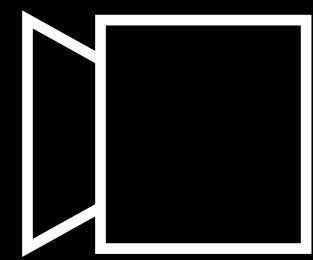


## Metrics

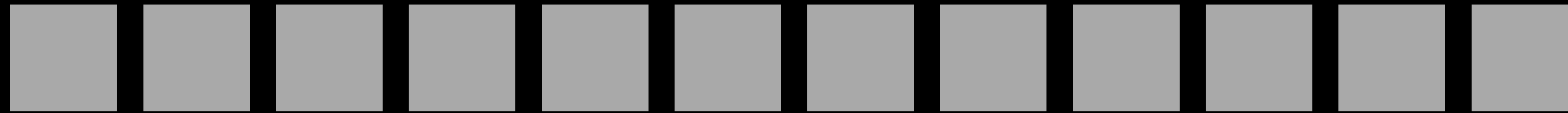
Precision

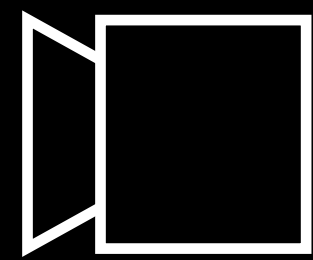
Recall

F-3

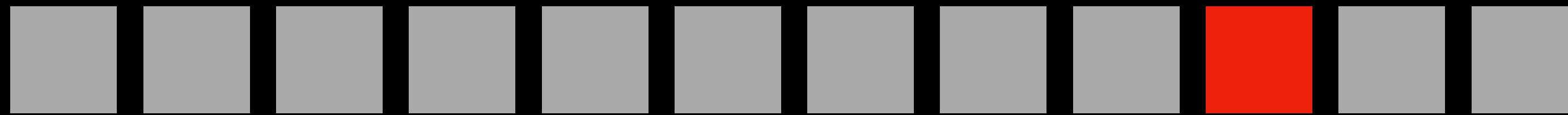


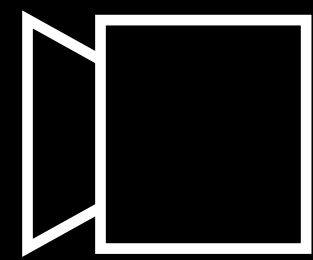
Camera



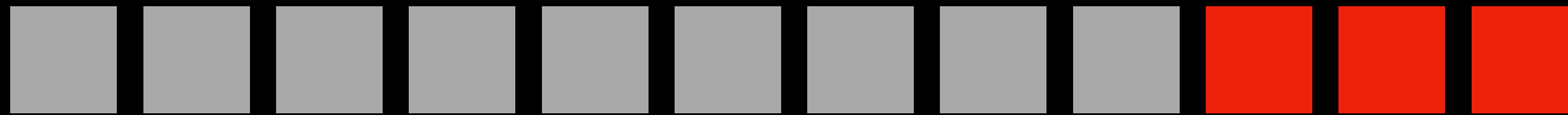


Camera



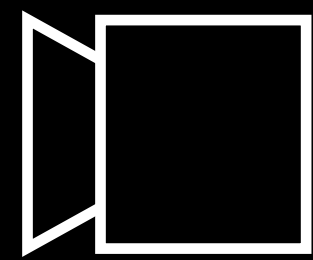


Camera

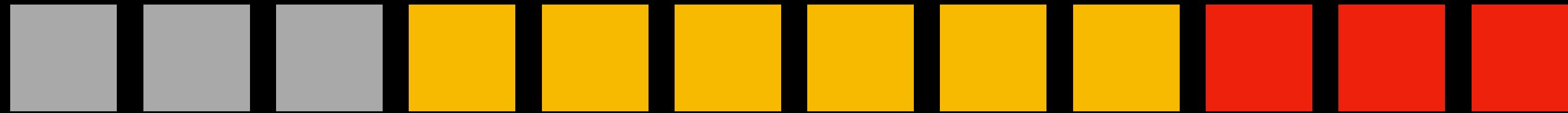


failure



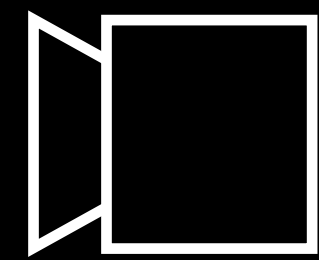


Camera

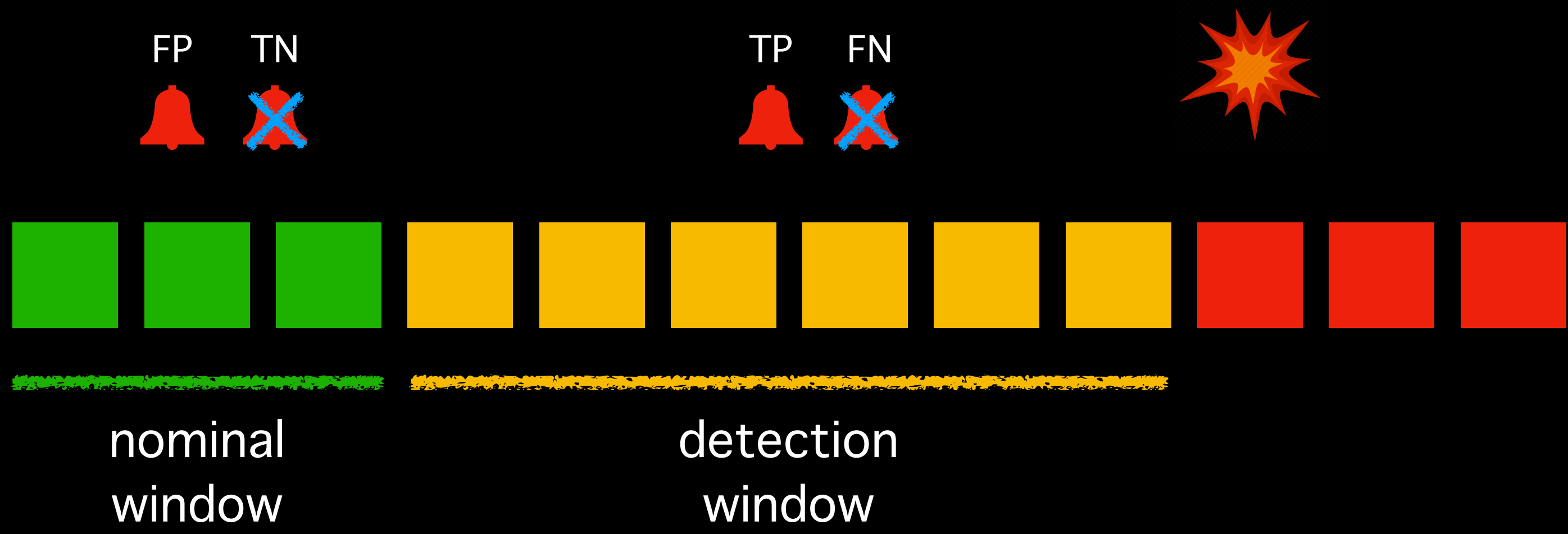


detection  
window

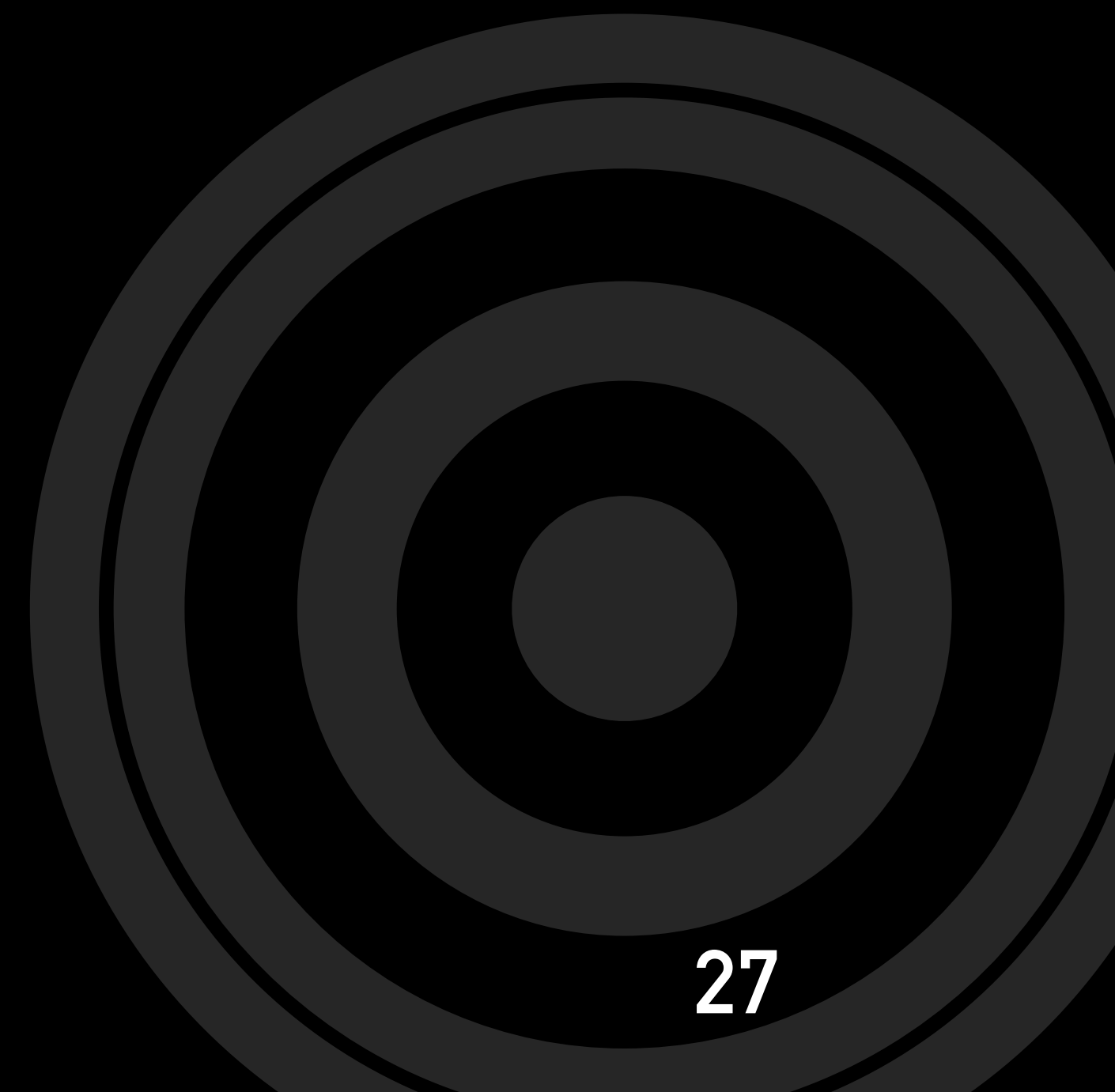




Camera



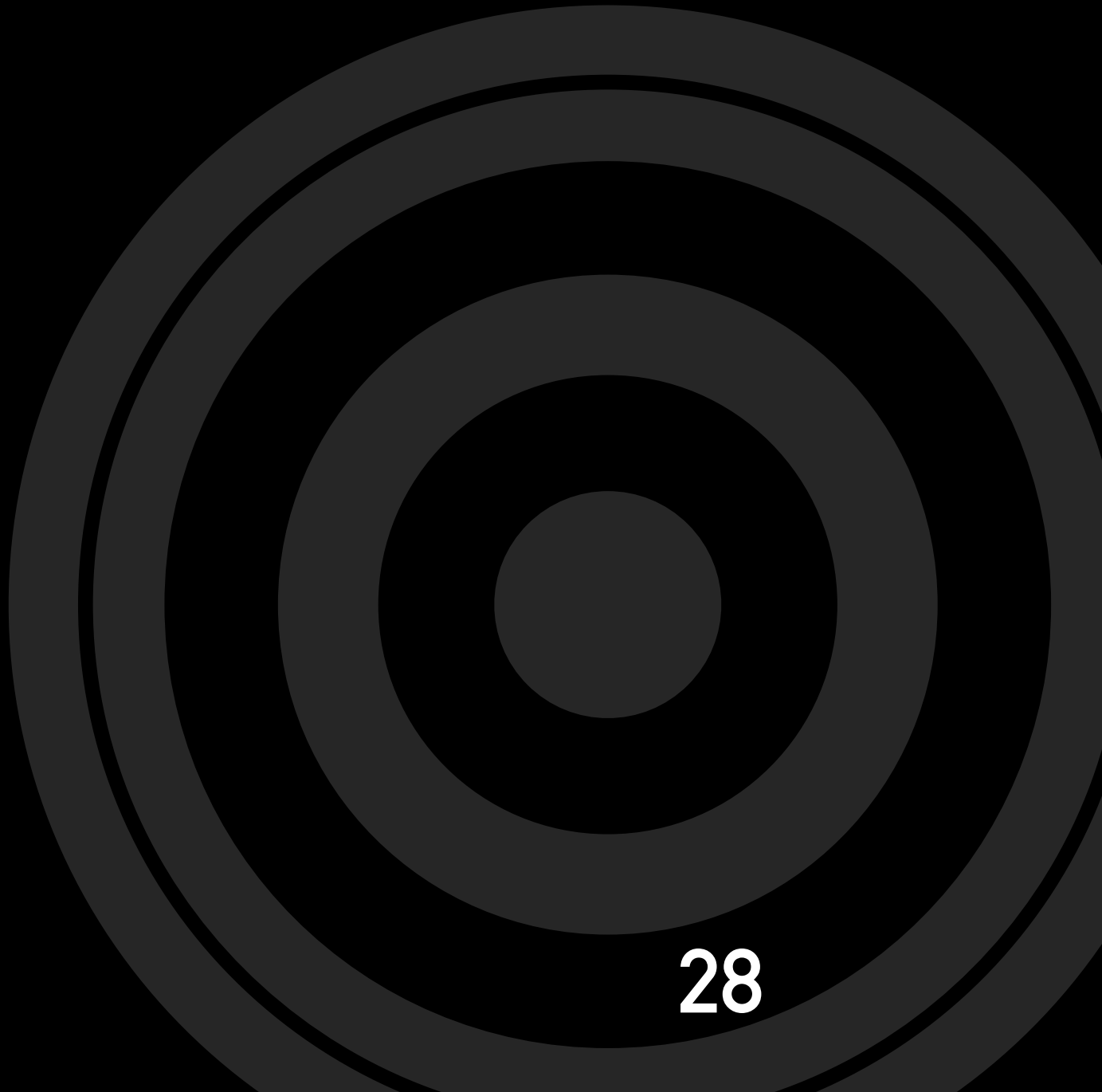
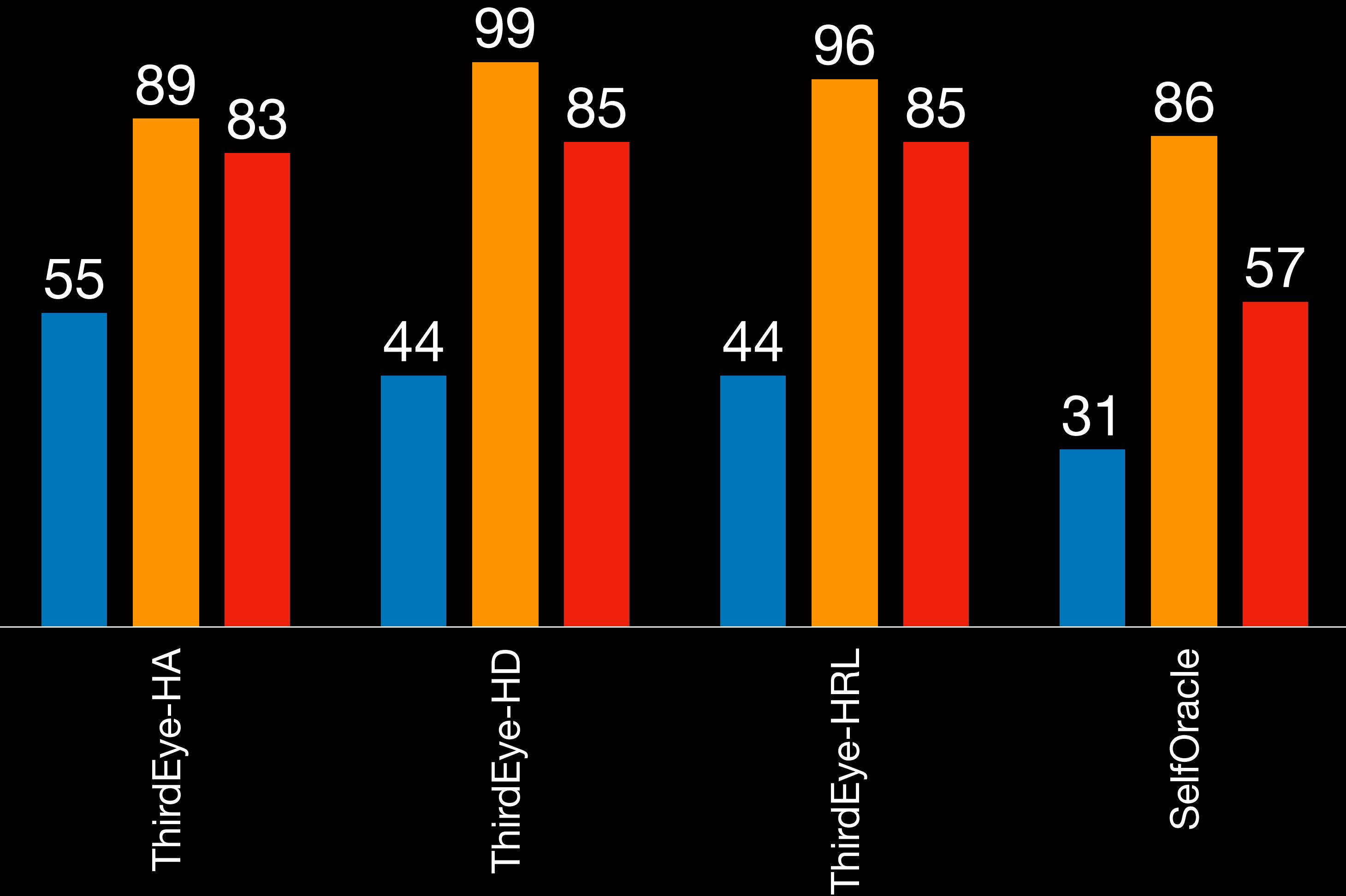
Is our approach **effective**  
in **predicting** misbehaviours ?

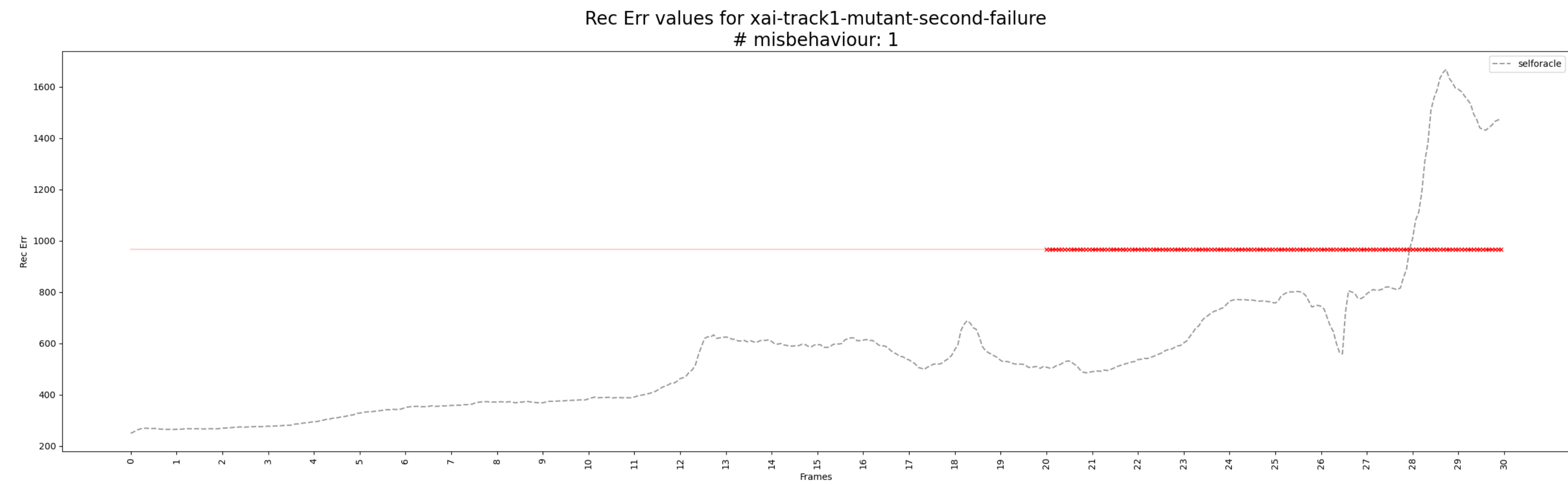
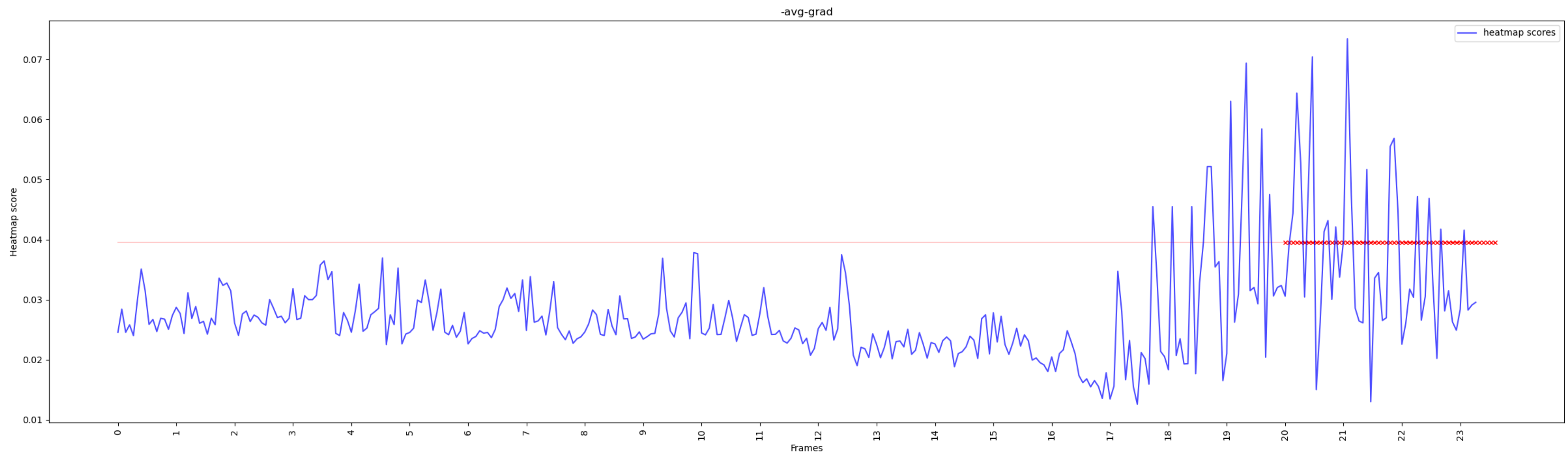


Is our approach **effective**  
in **predicting** misbehaviours



Precision  
Recall  
F-3  
 $\alpha=0.05$





*Can we **predict**  
**unexpected** conditions*



*ThirdEye enables **fast**  
and **accurate** online  
misbehaviour prediction*



Can we **predict**  
**unexpected** conditions ?

ThirdEye enables **fast**  
and **accurate** online  
misbehaviour prediction



## ThirdEye: Attention Maps for Safe Autonomous Driving Systems

### Code

<https://github.com/tsigalko18/ase22>

### Simulator

[https://github.com/tsigalko18/  
self-driving-car-sim/tree/  
USI\\_v1.0.0](https://github.com/tsigalko18/self-driving-car-sim/tree/USI_v1.0.0)