

Netherlands Forensic Institute  
*Ministry of Security and Justice*

## Forensic Intelligence workshop

prof. dr. ing. Zeno Geradts  
Senior forensic scientist /  
Special Chair Forensic Data Science  
Digital Technology and Biometrics /  
University of Amsterdam

DFRWS Oslo 2019



# COST Project DigForAsp

**DigForAsp (Digital forensics: evidence analysis via intelligent systems and practices)** – CA17124 is funded by the European Cooperation in Science and Technology (COST). DigForAsp activities were launched on 10th September 2018 for 4 years.



Funded by the Horizon 2020 Framework Programme of the European Union



## Outline

- Introduction
- Deep learning and neural networks
- Examples deepfakes
- Issues
- Outlook and conclusion





# Netherlands Forensic Institute





# University of Amsterdam

## Chair Forensic Data Science



- store and process
- understand and decide
- analyse and model
- Report and visualize
- Higher efficiency
- Data-intensive
- Evidential strength big data



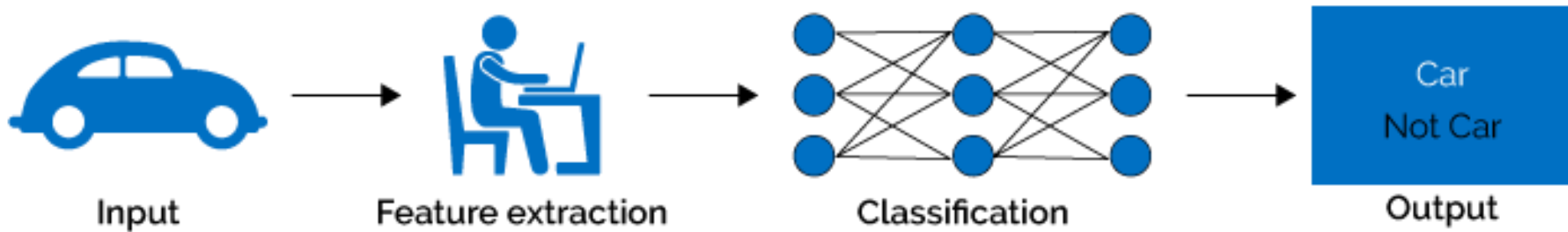
# Third Hype in history in AI



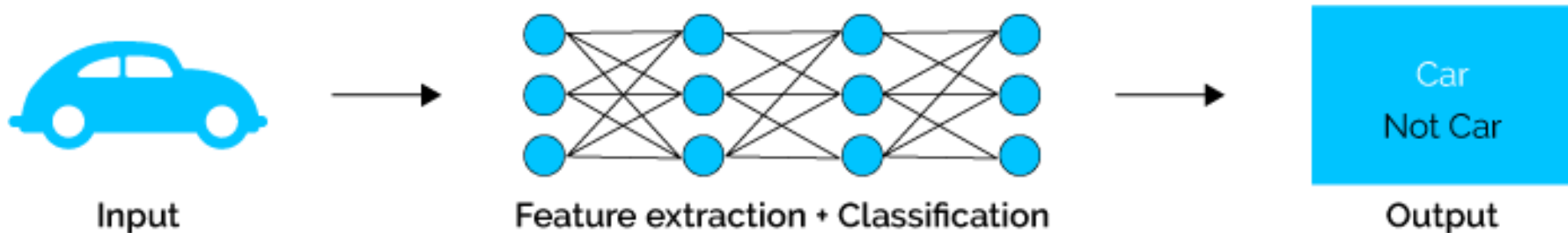


# Machine learning vs deep learning

## Machine Learning



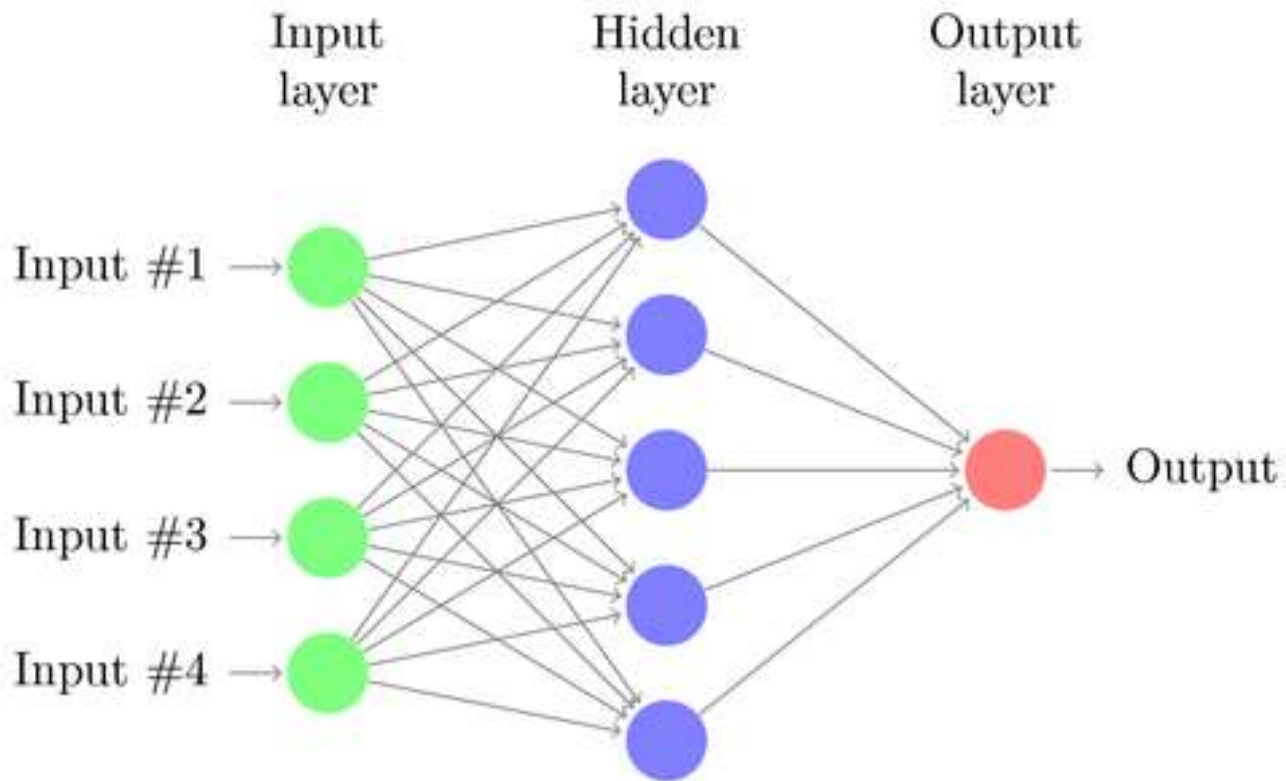
## Deep Learning







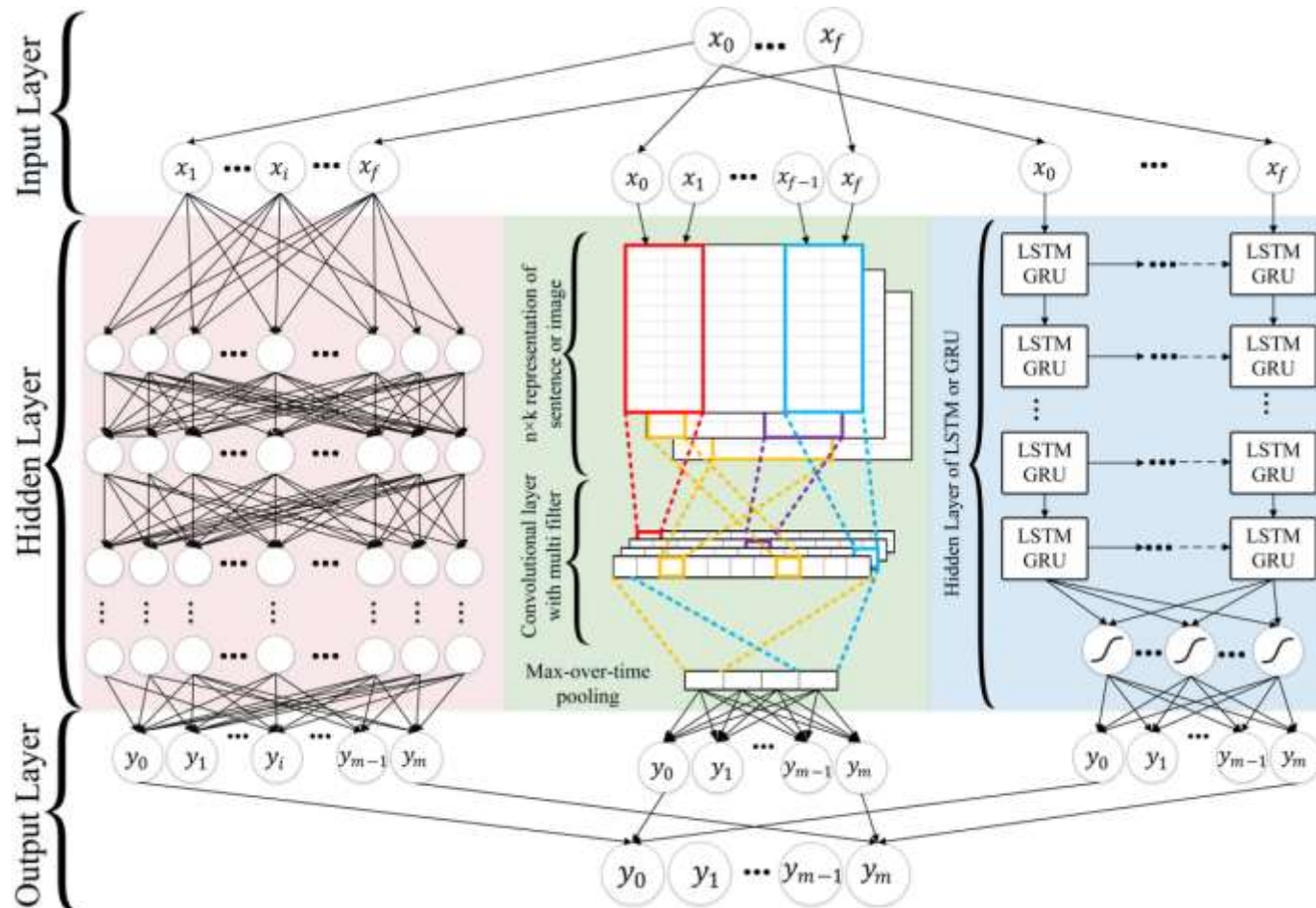
# neural network







# Neural network multilayer





# Calculation speed with Digital Evidence





# Digital Evidence





Extract data



Make data  
readable



Organize data



Interpret data





# Make data readable



## Challenge: many formats, old & new, non-standard

- Tool and library development

- Reverse engineering

Discover the technological principles of a system (e.g. software or communication protocol) through analysis of its function and operation

```
000025c0 0e 4a 5b fc 74 d6 21 a2 fb d3 5d bf 59 45 11 9a |.J[.t.!...].YE..|
000025d0 fd d6 00 28 d6 6f a1 b0 60 59 6c ce c6 d0 4c 5c |...(.o..'Yl...L\|
000025e0 61 10 8d cf 95 41 c1 3e b3 f3 62 ff 1b b0 fc dc |a....A.>..b.....|
000025f0 ea 5b fb 07 95 27 28 59 9a 05 e0 06 27 7b 2a 59 |.[...'(Y....'{*Y|
00002600 0e 43 72 1b ce 4b 1f 59 e2 ce d9 f3 86 34 5e f9 |.Cr..K.Y.....4^.|
00002610 38 d1 4a 0f 06 2e 70 66 c9 49 01 00 7b ca 93 c2 |8.J...pf.I.{...|
00002620 6d 70 02 ab b6 78 90 e1 5b ca 1c 14 29 13 77 93 |mp...x..[...].w.|
00002630 9f 29 a4 d1 1f 1f 3f 20 69 29 c4 ae fd c3 01 bf |.)....? i).....|
00002640 76 c4 bd a8 cc 99 0b e3 93 74 82 b8 1e cc 2e da |v.....t.....|
00002650 64 eb 74 64 5c 6c d7 91 78 5a 58 5b 59 c5 9a 82 |d.td\l..xZX[Y...|
00002660 4d e0 2c 58 1b 5c 83 c7 7e 98 3e 37 b2 93 99 90 |M..X.\...->7....|
00002670 fd 00 e0 3a 8e 4f 13 e5 1f 23 bb b5 f8 b0 a3 85 |...:0...#.....|
00002680 86 74 b9 1b 18 b7 5f 03 4b a1 6a c5 7c c4 46 1e |.t.....K.j|.F..|
00002690 6b 09 51 77 6b 3b 0d 9c 17 36 31 71 07 f4 9a bb |k.Qwk;...61q....|
```



## Trace Recovery & Analysis

Trace-analysis is the expertise to conserve, detect, repair, undelete, decrypt, find, structure and interpret data and traces on any case related digital medium.



e	4a	5b	fc	74	d6	21	a2	fb	d3	5d	bf	59	45	11	9a	.J[.t.!...].YE..	
d	d6	00	28	d6	6f	a1	b0	60	59	6c	ce	c6	d0	4c	5c	...(.o..`Yl...L\	
1	10	8d	cf	95	41	c1	3e	b3	f3	62	ff	1b	b0	fc	dc	a....A.>..b.....	
a	5b	fb	07	95	27	28	59	9a	05	e0	06	27	7b	2a	59	.[...'(Y....'{*Y	
e	43	72	1b	ce	4b	1f	59	e2	ce	d9	f3	86	34	5e	f9	.Cr..K.Y.....4^.	
8	d1	4a	0f	06	2e	70	66	c9	49	01	00	7b	ca	93	c2	8.J...pf.I..{...	
d	70	02	ab	b6	78	90	e1	5b	ca	1c	14	29	13	77	93	mp...x..[...].w.	
f	29	a4	d1	1f	1f	3f	20	69	29	c4	ae	fd	c3	01	bf	.)....? i).....	
6	c4	bd	a8	cc	99	0b	e3	93	74	82	b8	1e	cc	2e	da	v.....t.....	
4	eb	74	64	5c	6c	d7	91	78	5a	58	5b	59	c5	9a	82	d.td\l..xZX[Y...	
00002660	4d	e0	2c	58	1b	5c	83	c7	7e	98	3e	37	b2	93	99	90	M.,X.\..~.>7....
00002670	fd	00	e0	3a	8e	4f	13	e5	1f	23	bb	b5	f8	b0	a3	85	...:0...#.....
00002680	86	74	b9	1b	18	b7	5f	03	4b	a1	6a	c5	7c	c4	46	1e	.t...._K.j. .F.
00002690	6b	09	51	77	6b	3b	0d	9c	17	36	31	71	07	f4	9a	bb	k.Qwk;...61q....









### THE IoT PLATFORM OPPORTUNITY

# The Internet of Things (IoT) has a potential economic impact of 2.7-6.2 trillion USD until 2025

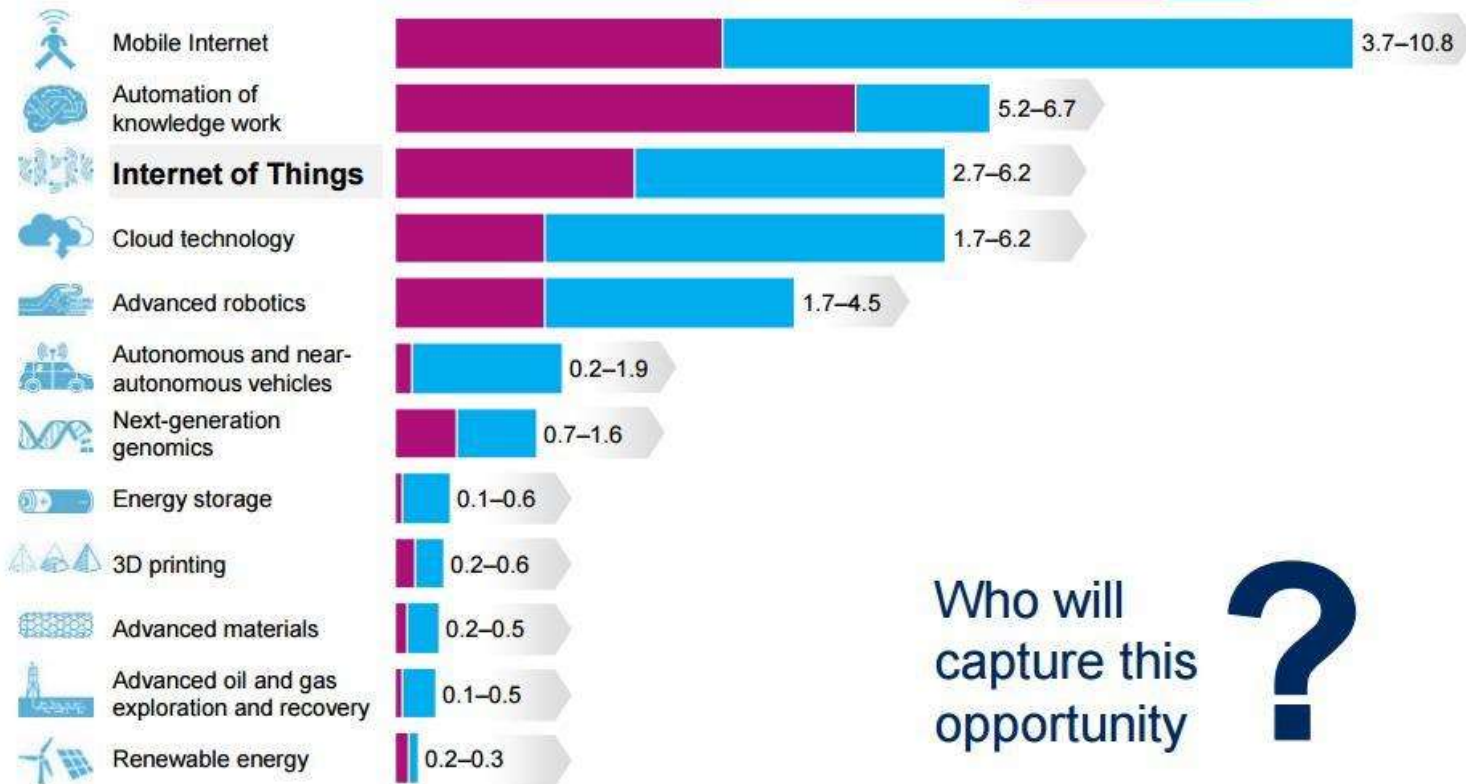
\$ trillion, annual

Range of sized potential economic impacts

Low High

Impact from other potential applications (not sized)

X-Y



Who will capture this opportunity ?

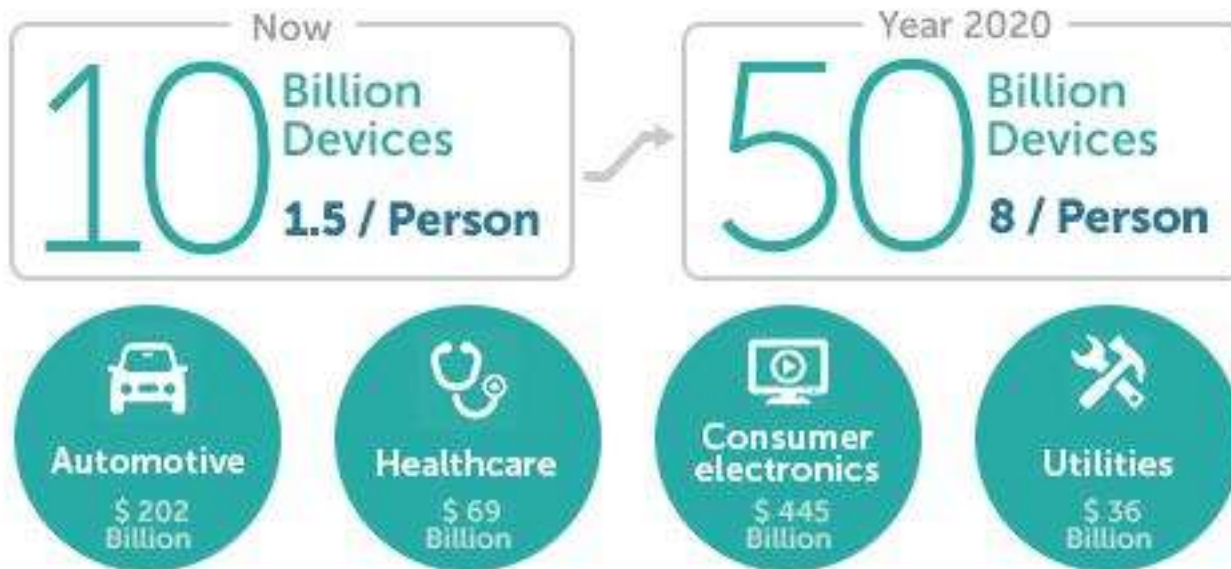
SOURCE: McKinsey Global Institute analysis

McKinsey & Company 3



# Internet of things 2020 Gartner

## IoT Predictions 2020





# 5G antenna and fiber boom





# Big Data issues





MIND & BODY, RESEARCH, TECHNOLOGY & ENGINEERING

## Everything big data claims to know about you could be wrong

By [Yasmin Anwar](#), Media Relations | JUNE 18, 2018



When it comes to understanding what makes people tick — and get sick — medical science has long assumed that the bigger the sample of human subjects, the better. But new research led by UC Berkeley suggests this big-data approach may be wildly off the mark.

That's largely because emotions, behavior and physiology vary markedly from one person to



TO



RE





#T&W14

# HOW NOT TO LET THE INTERNET KNOW YOU'RE PREGNANT: A 9-MONTH INFRASTRUCTURAL STUDY



NO SOCIAL MEDIA

PAY FOR EVERYTHING IN CASH

SHIP EVERYTHING PO. BOX

## The Server

"...BUT WHAT'S ON YOUR MAIL..."  
DON'T REALIZE THE DR. WASN'T PRIVATE...  
ALL DATA USEFUL TO STORE, TRACE & FOR YOU, AND FOR MARKETING!

The right for a transaction to be just a transaction

WITHDRAWING LOTS OF MONEY AT ONCE AS RED FLAGS

NOT A FREE MARKET, BUT COERCION

IN SCRIPTING THE HISTORY AS AMORAL...  
TOOLS USED TO CONNECT ALSO USED TO TRACK.

# SQUEAK DOLPHIN TO NORMCORE ANXIETY IN BIG DATA ERA.

What will the that reality of big data feel like?

BLENDING IN VIA COMFORTABLE SAMENESS.  
TRENDY WARE ALSO BECOMING IMPOSSIBLE

CANNOT PROCESS ALL THIS DATA  
THE DATA ITSELF IS NOT ENOUGH



The anxiety of being surveilled... but imagine what it's like to have all this data... and it not mean -ANYTHING-

# BIG DATA TO GROUND DATA



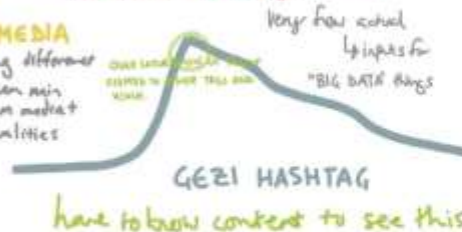
## Security/Privacy

BIG DATA PLUS SMALL, LOCAL, THICK DATA.

## CENSORSHIP ROUTED AROUND

#vizithink by @willow100  
VIZ.BLOOCYB.ORG

PENGUIN MEDIA  
seeing difference between main stream media + actualities  
SUBJECTS OUT OF CONTEXT, WON'T BE SCRAPPED



# TMI:

# BIG DATA IS ABOUT PEOPLE.

# THEORIZING BIG DATA



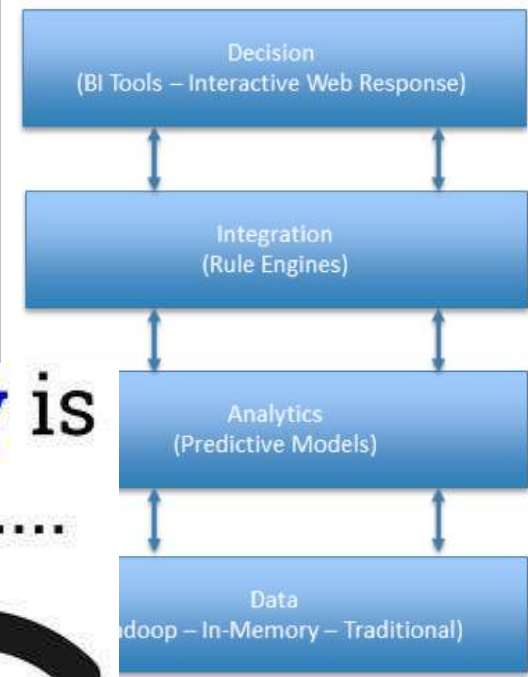


The good news : many examples where it works well credit card fraud detection and casework  
VISA states they save billions of euros a year

The **Fear** of **Technology** is **STRONG** with this one....



RTBDA for Predictive Analytics







# Big Data at NFI

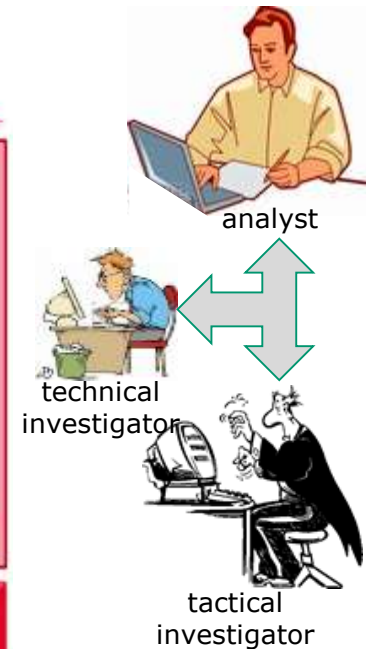
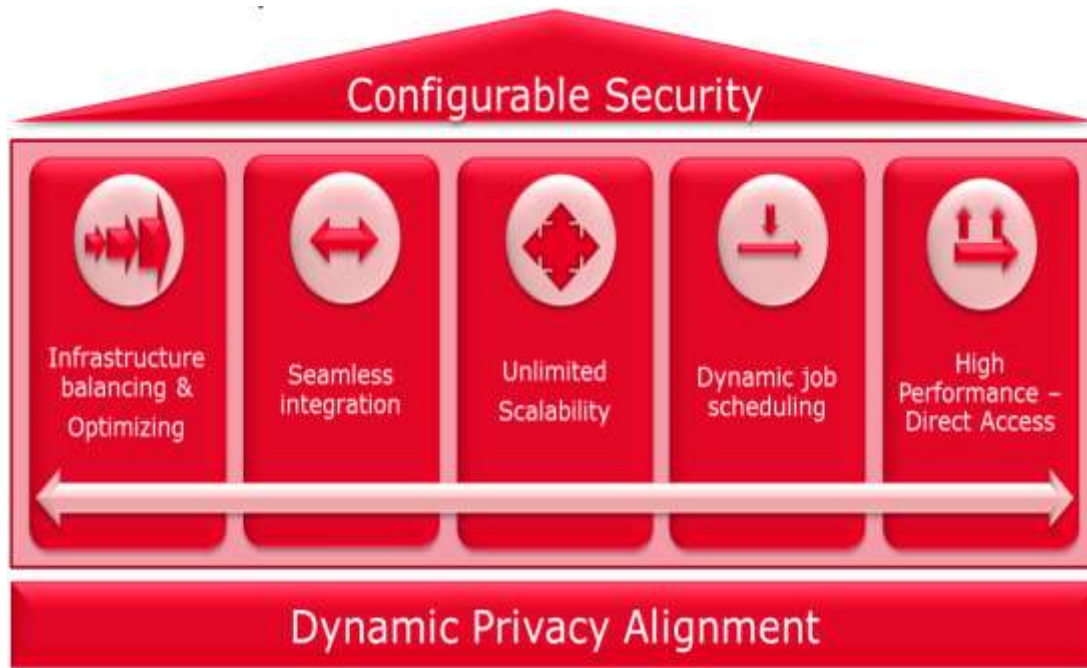
- Text Mining
- Data Profiling
- Financial Data Analysis
- Social Network Analysis
- video and images
- using big data analysis in forensic science





# Future of digital investigation: HANSKEN

CONFISCATION > SECURE > ENABLE ACCESS / ENRICH > REPORT > ANALYSE



Some hours (1Tb/20 min) – direct results at start



# Evolution forensic analysis – automation, speed & coverage

manual  
import and  
manual  
processing

**Conventional:** throughput months

50%

50%

manual  
import and  
automated  
processing

**XIRAF:** throughput weeks

70%

30%

automated  
import and  
automated  
massive-  
parallel  
processing

**HANSKEN:** throughput hours

85%

15%





# Many new techniques



Puppet

**Nagios**



Digital Investigation 11 (2014) S54-S62



Contents lists available at ScienceDirect

# Digital Investigation

journal homepage: [www.elsevier.com/locate/diin](http://www.elsevier.com/locate/diin)



## Digital Forensics as a Service: A game changer



R.B. van Baar\*, H.M.A. van Beek, E.J. van Eijk

*Netherlands Forensics Institute, Laan van Ypenburg 6, 2497 GB The Hague, The Netherlands*

### A B S T R A C T

**Keywords:**  
Digital forensics  
DFaaS

How is it that digital investigators are always busy and still never have enough time to actually dig deep into digital evidence? In this paper we will explore the current implementation of the digital forensic process and analyze factors that impact the efficiency of

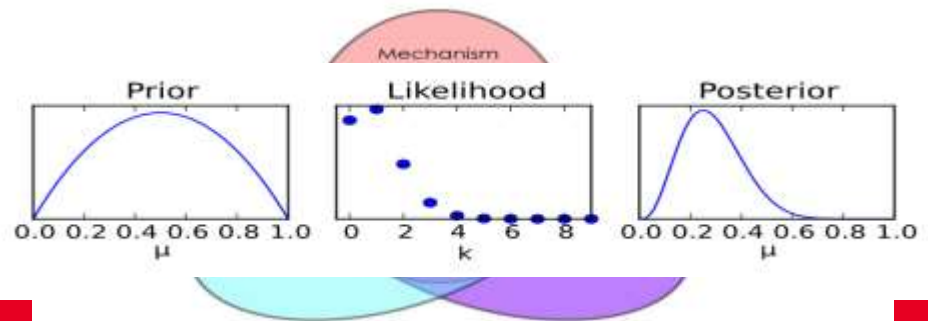
is processed and of the manner in which the traces collected by this processing is analyzed

**Related work**



# Examples hypotheses in digital forensic science

- has the computer been hacked or not ?
- has the email been send or not ?
- has the USB been plugged in or not ?
- was the phone in this location or at the location presented by the defence ?
- has the child pornography been send by the computer of the suspect or not ?
- is the child porn photographed with this camera or another camera ?







# Lawful internet interception







# IMSI catcher : privacy by design ?

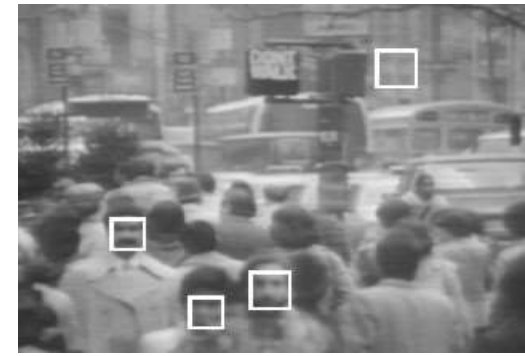




## Challenge: data is not self-explaining

Add models and analysis to support interpretation

- Scenario analysis
- Timeline analysis
- Geographical models: e.g. location of cell phones
- Analysis of images / video / audio
  - Size
  - Speed
  - Face recognition
  - Speech recognition
- Author recognition





## About Forensic Big Data Analysis

- Our data come from confiscated phones, hard drives, licence plate cameras, telephone providers, and so on...





## What if...

- An ATM machine is blown up
- A prepaid cell phone is found on the scene
- The police have their eyes on a suspect
- **Research question: is the suspect the user of the prepaid phone?**





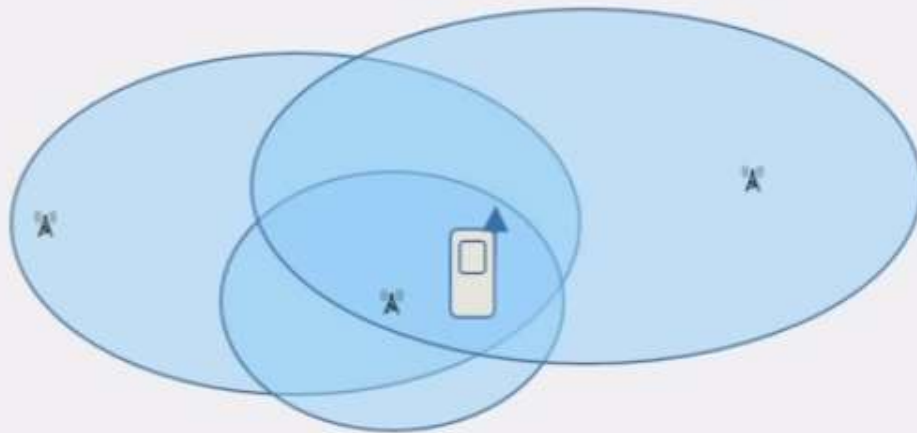


## What information do we have?

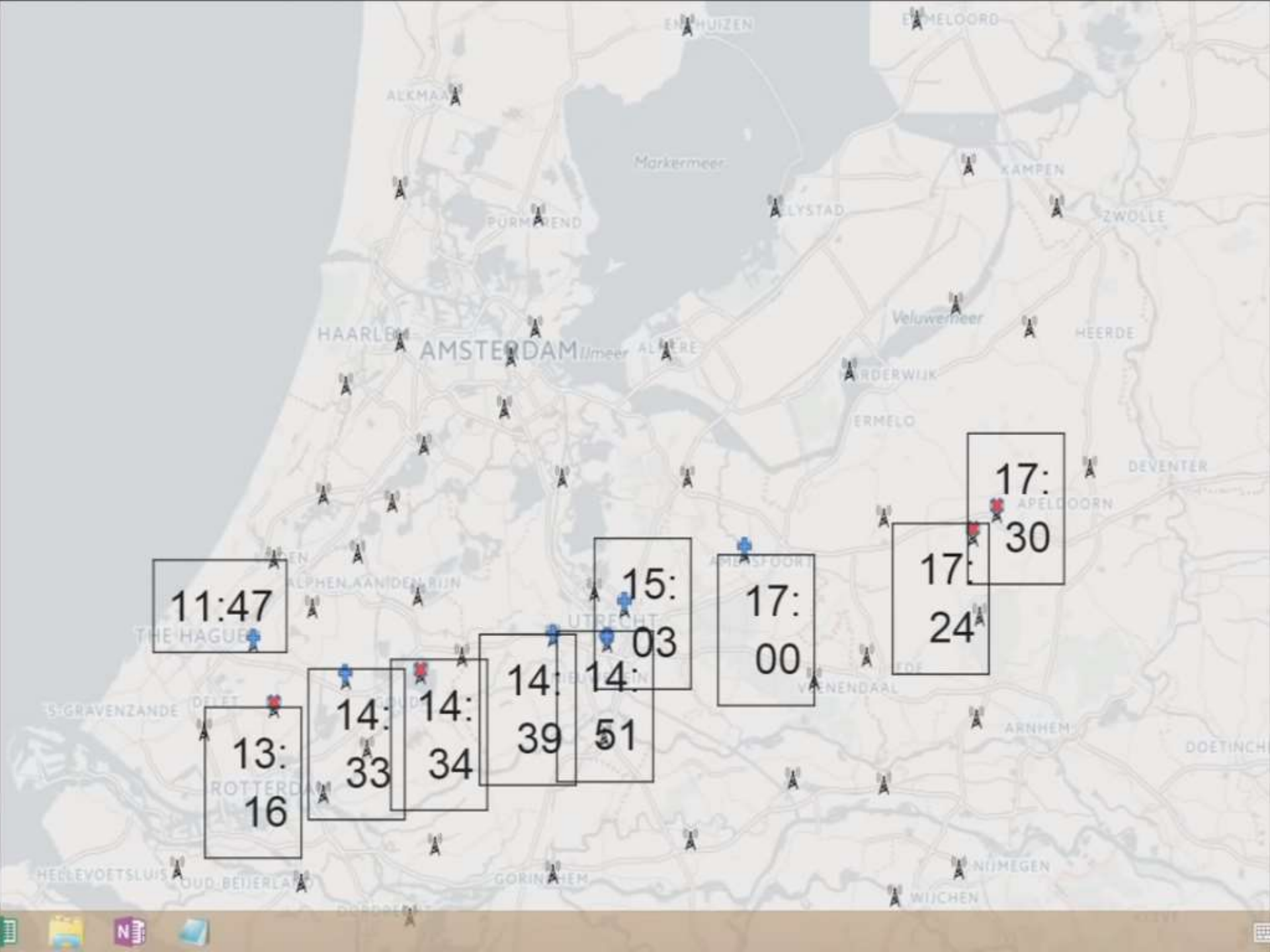
- You know the phone number of the prepaid phone and that of the suspect's private phone.
- The telephone provider provides the police with usage data for both phones.
- Every time a phone connects to a cell tower, you know when it happened.
- You know the location of each cell tower.



## Problem 1: cell tower location data are not precise and depend on...



- Theoretical range: 35km
- Direction of transmission
- Distance
- Obstacles (tall buildings)
- Weather conditions
- Network load



11:47

13:16

14:33

14:34

14:39

14:51

15:03

17:00

17:24

17:30



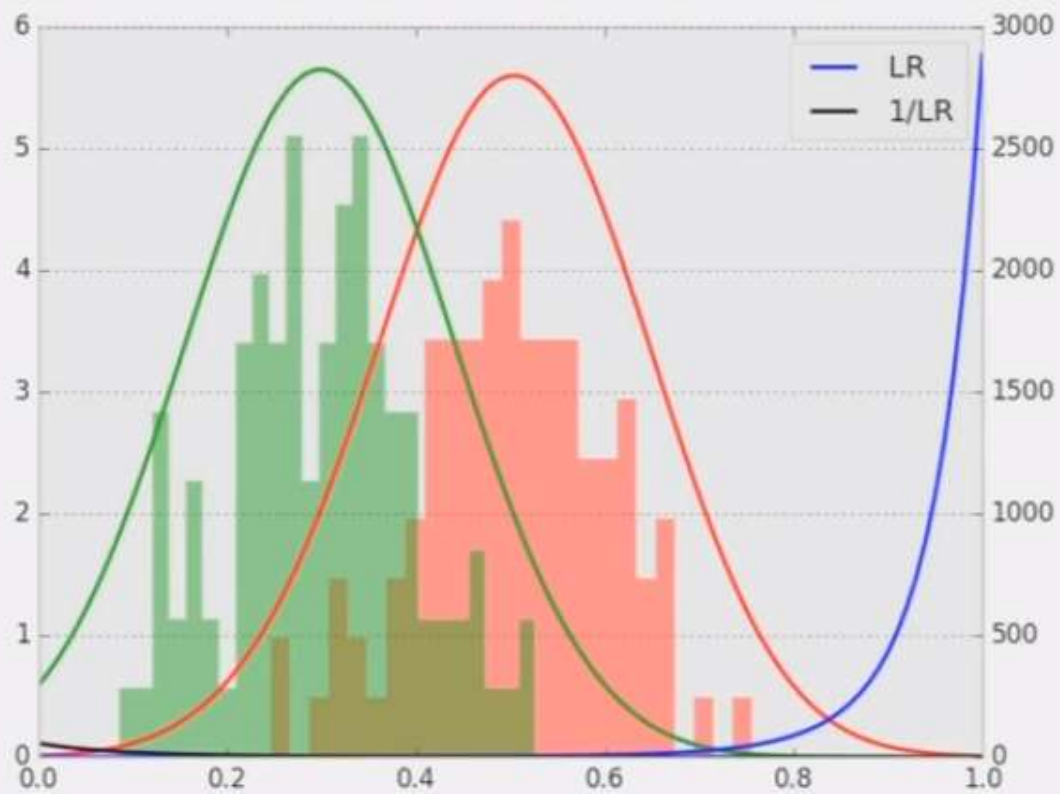


## To summarize...

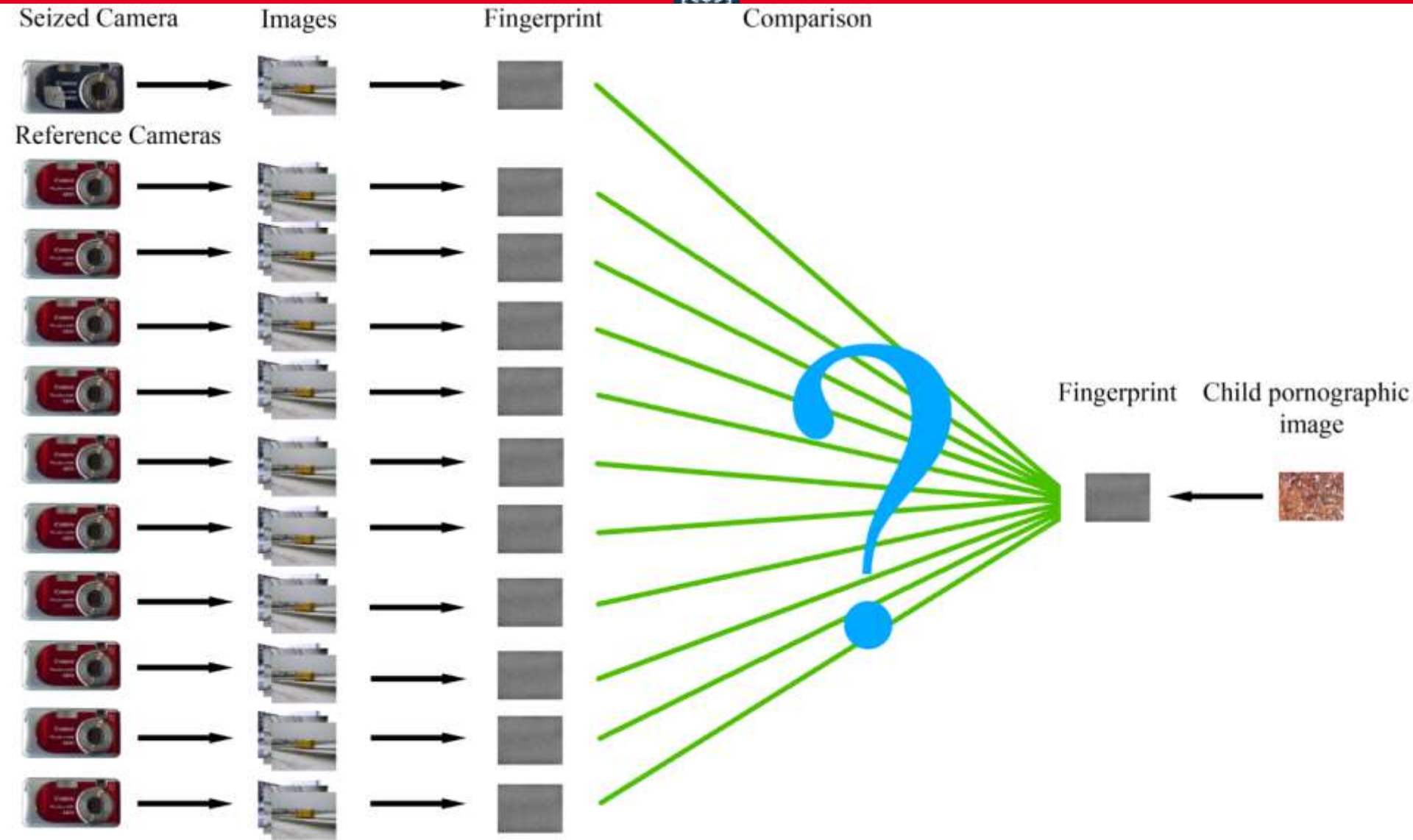
- We want to know if the suspect is the user of a prepaid phone that can be linked to a crime.
- We know when and where the prepaid phone was used.
- We know when and where the suspect's phone was used.
- But our data are sparse and imprecise...



# Likelihood Ratio



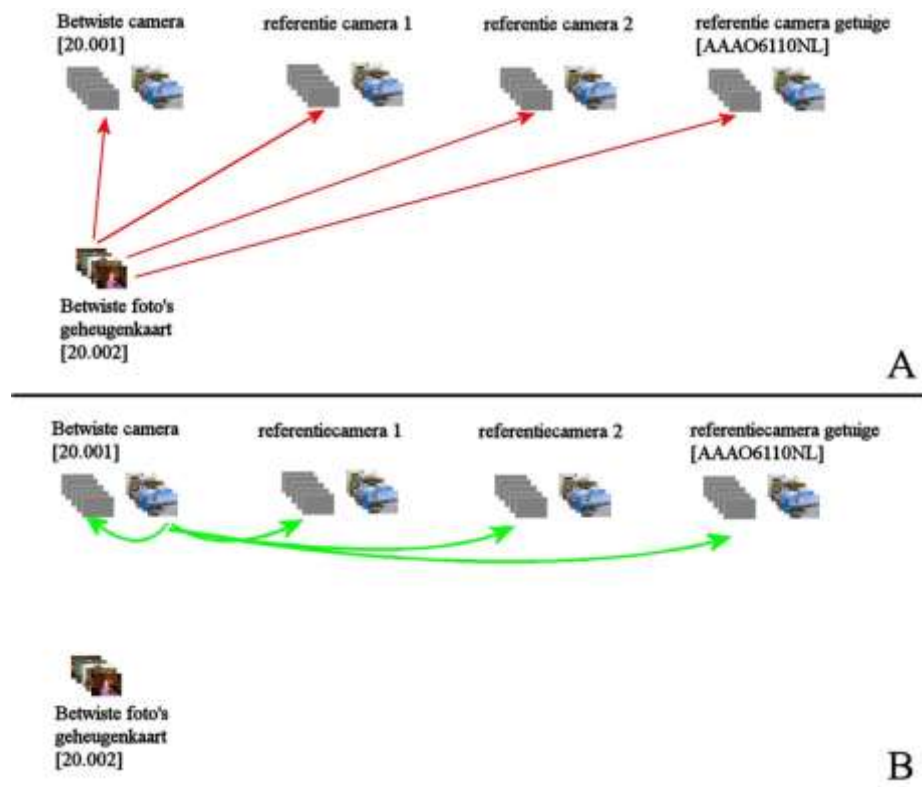








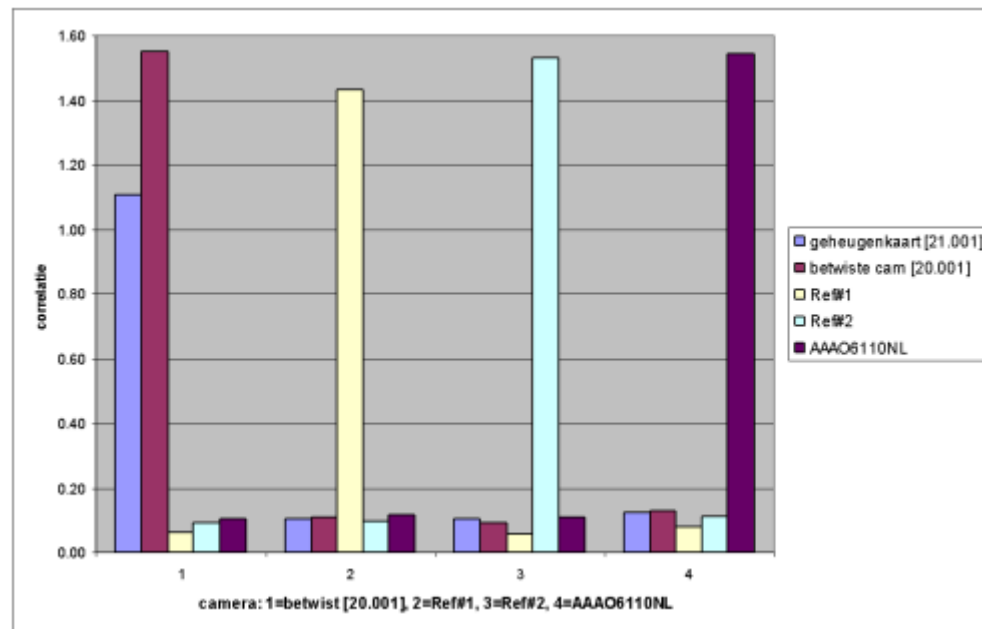
# Casework links





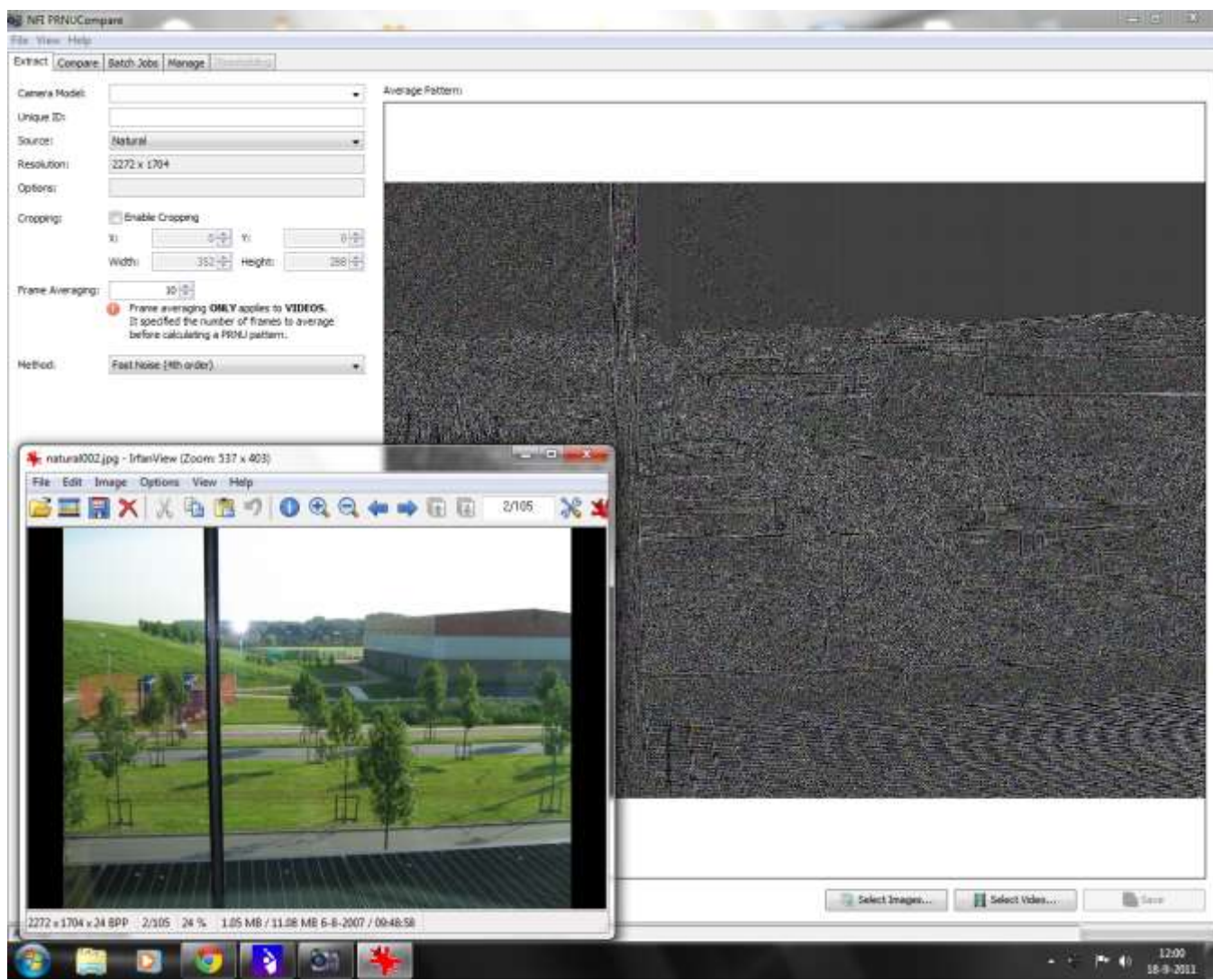
# Casework

- Example where it worked





# PRNU Compare





# Bayesian

Question: were the images made with the seized camera?

Conclusion

The findings of the investigation are:

Equally likely

Somewhat more likely

More likely

Much more likely

Very much more likely

if H1 is true, than if H2 is true.

The findings are very much more likely if the Seized Camera took the child pornographic image, than if another camera took the image.

Camera	...	...	...	Sum
<b>Comparing: Foto in kwestie (Onbekend)</b>				
Verdachte Camera reference (Canon PowerShot)	..	..	...	0.131330
Camera 1 reference (Canon PowerShot)	..	..	...	0.008054
Camera 4 reference (Canon PowerShot)	..	..	...	0.007700
Camera 2 reference (Canon PowerShot)	..	..	...	0.007022
Camera 3 reference (Canon PowerShot)	..	..	...	0.006287

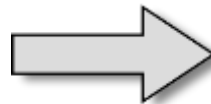




# Large Scale Camera Identification

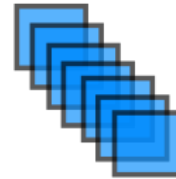
- Sorting photos by source
- Identify photos from the same source (camera)
- New valuable information and insight

unsorted photos

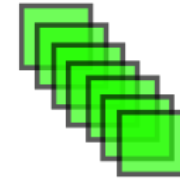


Panda

source 1



source 2



source 3



source 4





# Sorting Images by Source

**Scan** → Extract → Compare → Cluster → Explore



Scan

4320x3240

1024x768



Sorted by resolution and directory



# Sorting Images by Source

Scan → **Extract** → Compare → Cluster → Explore



Extract

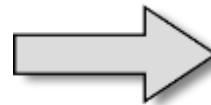
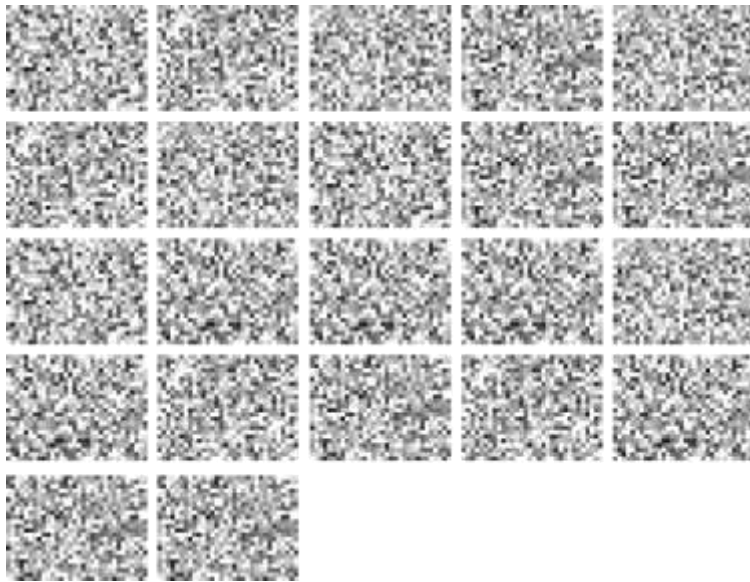


PRNU noise patterns (fingerprints)



# Sorting Images by Source

Scan → Extract → **Compare** → Cluster → Explore



Compare



Images compared to all images





# Sorting Images by Source

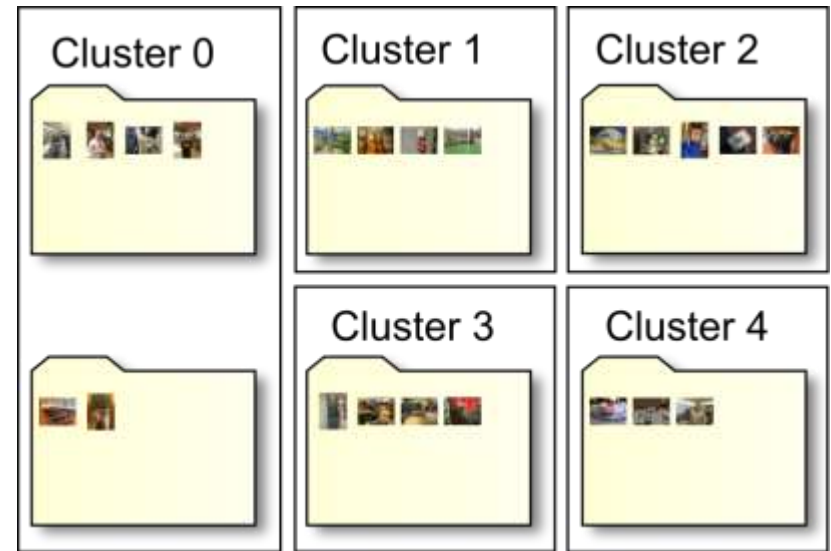
Scan → Extract → Compare → **Cluster** → Explore



threshold = 0.001



Cluster

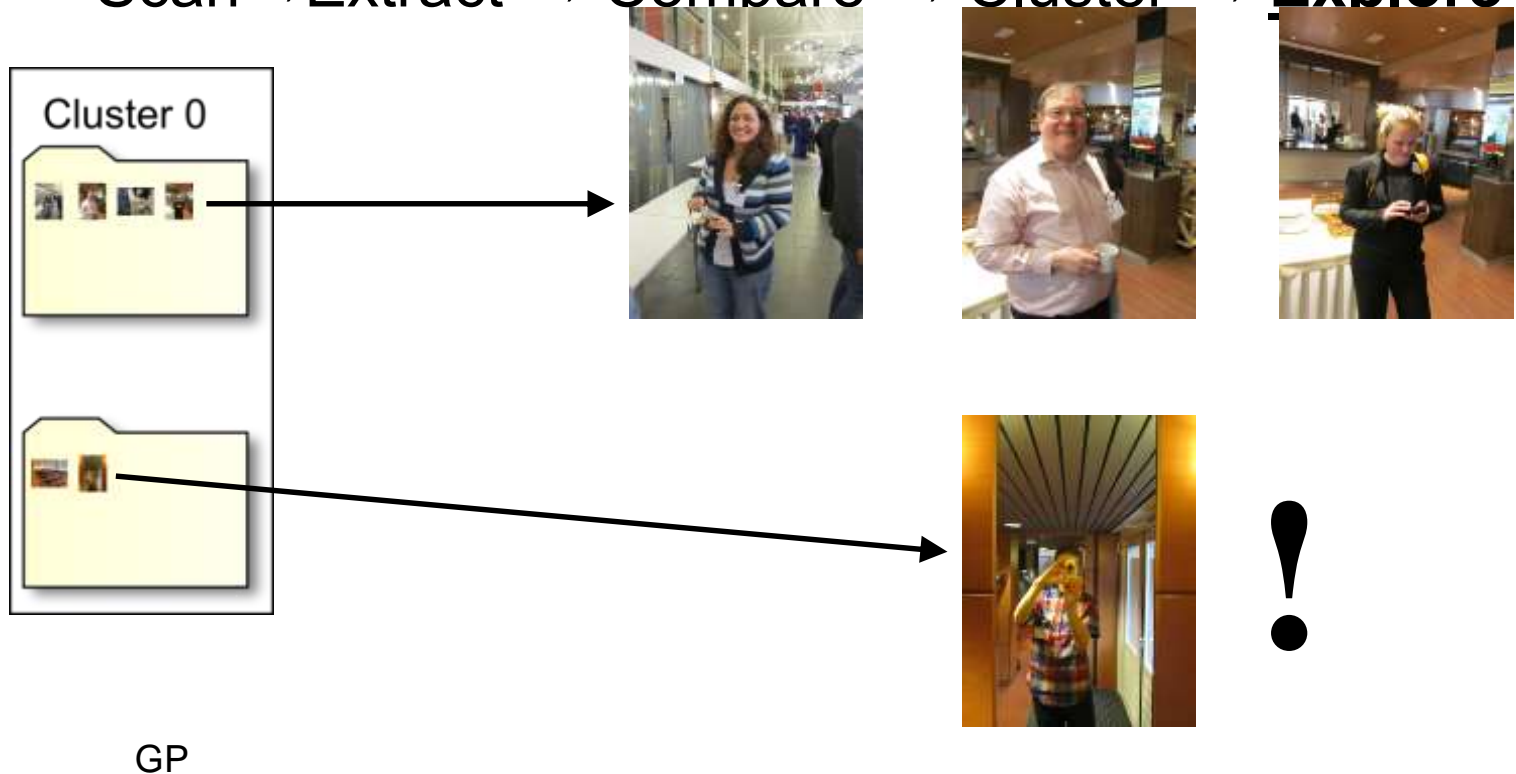


Images grouped by source



Sorting Images by Source also GPU / social networks also deep learning applied

Scan → Extract → Compare → Cluster → **Explore**





# Facial comparison



**NO**

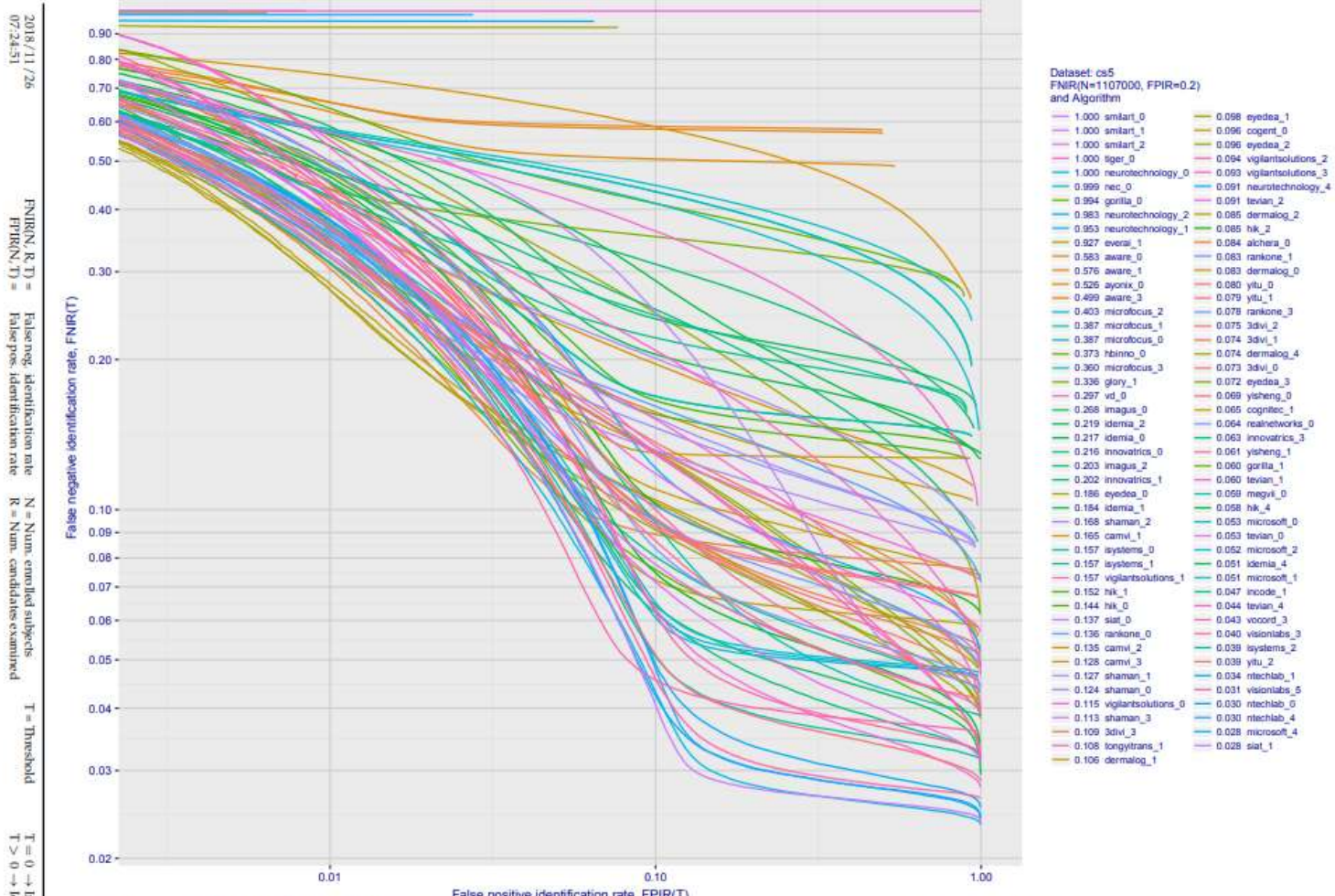


**YES**



# NIST test of faces in the wild

This publication is available free of charge from: <https://doi.org/10.6028/NISTIR.8238>



FRVT - FACE RECOGNITION VENDOR TEST - IDENTIFICATION

2018/11/26  
07:24:51

FNIR(N, T) =  
FPIR(N, T) =

False neg. identification rate  
False pos. identification rate

N = Num. enrolled subjects  
R = Num. candidates examined

T = Threshold

T = 0 → Invert  
T > 0 → Identif

Figure 99: [Wild Dataset] Identification miss rates vs. false positive rates. The figure shows accuracy of algorithms on wild images searched against wild images of familiar individuals enrolled into gallery. On the horizontal axis is the False Positive Identification Rate (FPIR(N, T)) with N = 1107000 as a function of false positive identification rate (FPIR(N, T)).



# Other examples of deep learning

- manipulation detection
- face morphing / deepfakes
- court findings finding irregularities





# Detecting face morphing in video and documents

30 August 2018

Ilias Batskos

Andrea Macarulla Rodriguez

Marissa Koopman

Zeno Geradts

EAFS 2018 Lyon



# Contents

- Face Morphing
- Deepface
- Conclusions



# Definition and problem statement

- Morph: A novel photo created by blending the photos of two different individuals



Face **a**



Morph of  
**a** and **b**



Face **b**

- Problem: 2 individuals, one being the criminal and the other the accomplice, can use the same travel document



## Security issues

Detection can be performed in 2 stages by humans, computers or both

- Issuing stage: Seemingly flawless morphs can be accepted as genuine photos, ~50% FAR when unaware and ~20% FAR when aware (Issuing officer) [1]
- ABC stage: Morphs can bypass Automatic Border Control ( Face recognition systems),FaceVACS, FAR = 83.62% [2]

In real case scenarios of both stages, there are two photos that could be compared:

- Issuing stage : Photo from previous ID/Passport and presented new photo
- ABC stage : Passport photo and probe photo



## Example 2

- Genuine or morph ?



Probe



e-Pass





## Example 3

- Genuine or morph ?



Probe



e-Pass



## Example 4

- Genuine or morph ?



Probe



e-Pass



## Example 5

- Genuine or morph ?



Probe



e-Pass



## Example 7

- Genuine or morph ?



Probe



e-Pass



## Problem/Vulnerability

The possibility to provide printed photographs creates a vulnerability in the system, which can be exploited by criminals with face morphing skills.

## Solution?

- Live photograph enrollment at an authorized facility
- Secure Police web application for enrollment
- Additional biometric features (fingerprint, iris) on the chip
- Or better detection







## State of the Art detection

- Texture based features using Local Binary Patterns (LBP), Binarised Statistical Image Features (BSIF), Image Gradient magnitude (IG) , Local Phase Quantitation (LPQ), blind/referenceless image spatial quality evaluator (BRISQUE)
- Double jpeg compression detection: Benford features , DCT coefficients of JPEG compressed face images
- Neural networks

### Examples



# SoA limitations

- Vulnerable to image processing and print & scan process
- Ghost artifacts, interpolation effects, morphing traces can be mitigated by image processing.(If the feature space of a detector is known it can be bypassed )
- Crucial pixel information is lost during print & scan, resulting to significantly increased errors of SoA detection methods



Ghost artifacts



Processed



Printed and Scanned



Digital photo

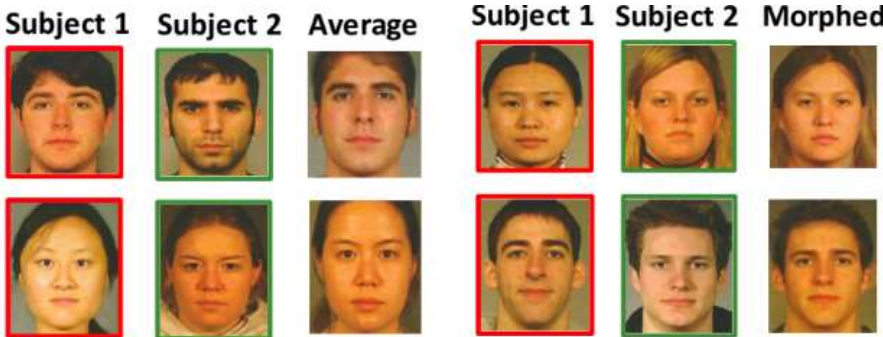


# Candidates selection

- Face encodings(128d) were extracted from all candidates.(Resnet34, Dlib)
- Using the Euclidean distance, a list of distance scores was calculated for each candidate.
- The person whose face encoding was closer to the candidate's encoding was selected for morphing



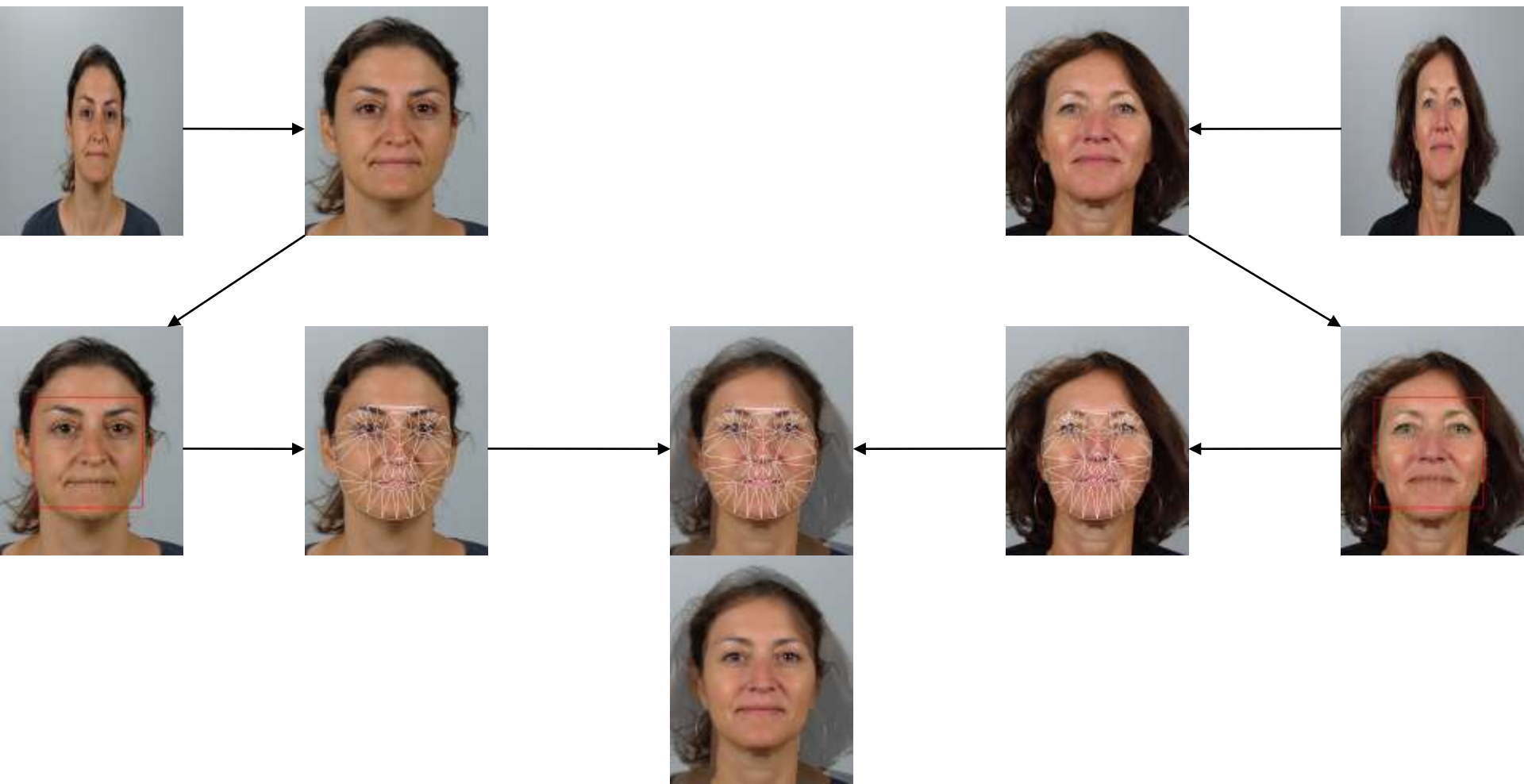
Automatically selected candidates according to similarity



Morphing candidates found in the literature



# Morphing pipeline





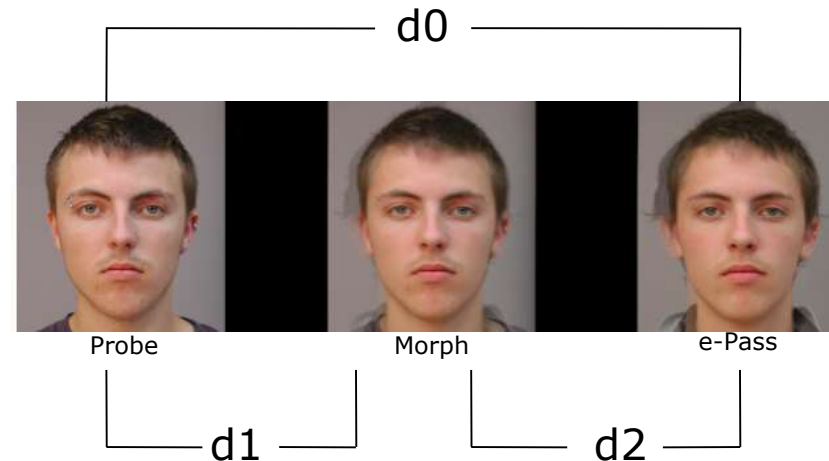
# Experimental method

## Criminal case

In a criminal scenario the photograph intended for e-Pass use already contains 50% of another individual, thus the morph will still contain 25% of that other individual.

After creating the morph which is intended as the e-pass photo the same process as before is followed

Note that the photo used for morphing was taken at a different time than the probe photo to simulate a real case scenario







## Experimental method

- Training set:  $H_0=56$ ,  $H_1=55$
- Testing set:  $H_0=163$ ,  $H_1=359$
- Total  $H_0=219$
- Total  $H_1=414$



# Experimental results

- True Positives (Genuine photos classified as genuine) = 146
- False Negative (Genuine photos classified as morphs) = 17
- Total Positives (Genuine photos) = 163
  
- True Negatives (Morphs classified as morphs) = 353
- False Positives (Morphs classified as genuine) = 6
- Total Negatives (Morphs) = 359
  
- True Positive Rate (Recall) = 0.8957055214723927
- True Negative Rate = 0.9832869080779945
- False Positive Rate (1-TNR) = 0.016713091922005572
- False Negative Rate (1-TPR) = 0.10429447852760736
  
- Precision = 0.9605263157894737
- False Discovery Rate (1- Precision) = 0.039473684210526314
- 
- Accuracy (proportion of true results) = 0.9559386973180076
- F1 score = 0.926984126984127



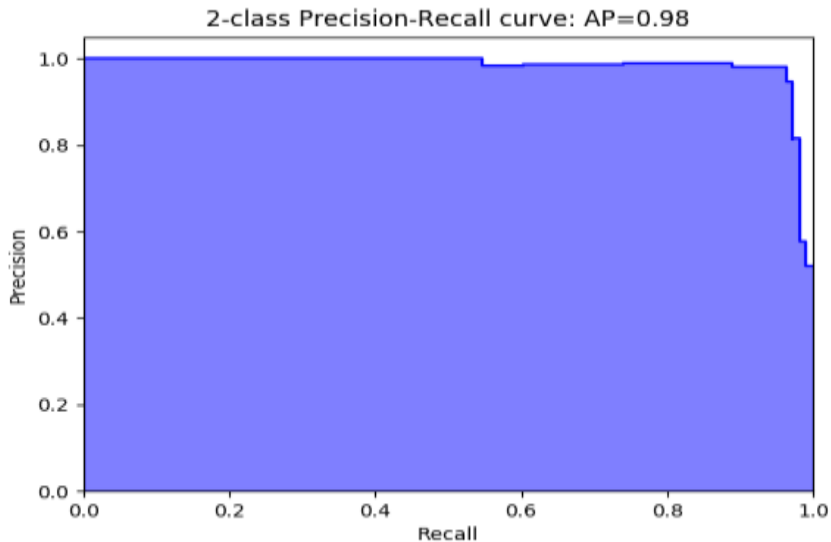
# Cross-validation

10-fold accuracy scores= [0.98245614, 1, 0.94736842, 0.96491228, 0.96428571, 0.92727273, 0.96363636, 0.98181818, 0.94545455, 0.98181818]

Accuracy: 0.97 (+/- 0.04)

Average precision-recall score: 0.98

auc= 0.989441195582

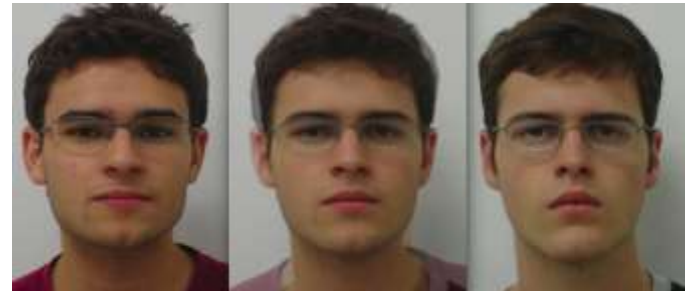


ROC curve



# Classification examples

False positives





# Classification examples

False negatives







# Classification examples

True negatives





## Limitations

- Sensitive to extreme photometric variations
- Sensitive to pose variations
- Sensitive to angle variations



# Improvements

- Cluster faces by gender and ethnic characteristics for better morphing candidates
- Ratios of distances between pairs of landmarks to create the face descriptor
- Additional landmarks to improve morph quality
- Second morphing, additional discriminating power(?)



## Conclusion

- The scores are unaffected by manual morphing(better quality than automatic) and Print and Scan process(?)
- Good experimental results indicate the effectiveness of the proposed method. The method could be implemented in parallel with SoA detection methods as an additional protection layer to counter cases of highly sophisticated and skilled criminals.

# Contents

- Introduction
  - Introduction to Deepfakes
  - Forensic relevance
  - Research goals
- Discrete Cosine Transformations (DCT)
  - Method
  - Results
  - Conclusions
- Convolutional Neural Networks (CNN)
  - Method
  - Results
  - Conclusions
- Photo Response Non Uniformity (PRNU) Analysis
  - Methods
  - Results
  - Conclusions
- Limitations
- Further research



# AI-Assisted Fake Porn Is Here and We're All Fucked

Someone used an algo \*\*\*\*\* ce  
of 'Wonder Woman' sta. ... porn  
video, and the implications are terrifying.

WEB / TECH / ARTIFICIAL INTELLIGENCE

## Gfycat starts removing fake AI porn GIFs from its platform

Gfycat is taking a firm stance against permitting 'deepfakes,' or AI-generated fake pornography

By Nick Statt | @nickstatt | Jan 31, 2018, 4:01pm EST

f t share

# Deepfakes: People are now swapping their friends' faces into porn

Deepfakes use facial recognition technology to superimpose faces on to porn stars.

By Kashmiri Gander  
January 29, 2018 15:30 GMT

TECH / ARTIFICIAL INTELLIGENCE / SEX

## Is it legal to swap someone's face into porn without consent?

Yes, no, maybe

By Megan Fackelmanesi | @Megan\_Nicole1 | Jan 30, 2018, 2:30pm EST



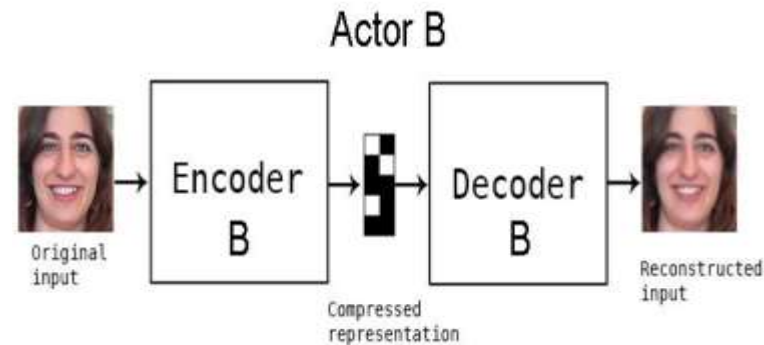
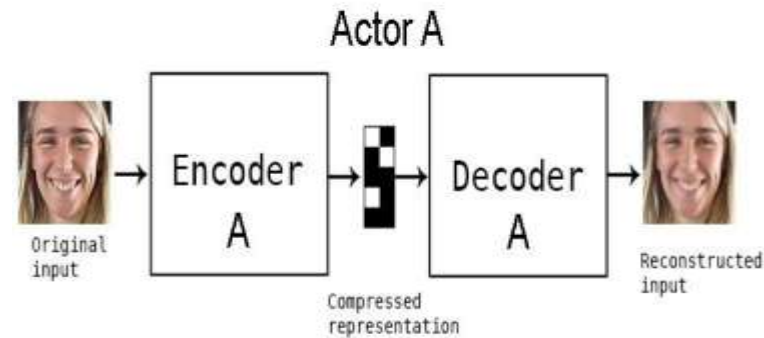
**DEEP**

**FAKES**

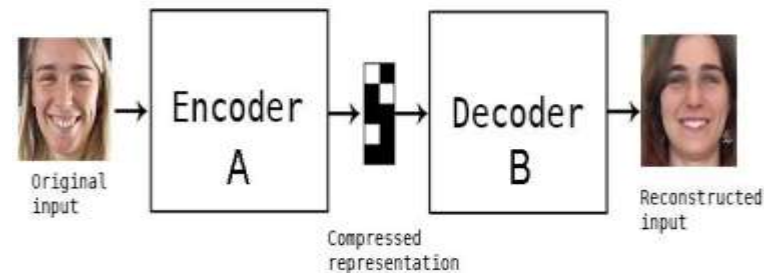
# How are Deepfakes made?

- Thousands of images of actor A and of actor B are needed for good results.
- Two autoencoders (A and B) are trained on these images.
- Autoencoder AB puts face of A on body of B

## Training



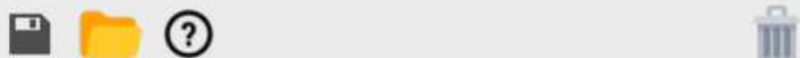
## Applied



# Forensic Relevance

- Authenticate video evidence
- Technology widely accessible on the internet
- Easy to use for everyone through GUIs
- Authentication of videos is also important for journalism, social media, etc.





VIDEO A	<input type="text"/>			
IMAGES A	<input type="text"/>			
FACES A	<input type="text"/>			
VIDEO B	<input type="text"/>			
IMAGES B	<input type="text"/>			
FACES B	<input type="text"/>			
MODEL	<input type="text"/>			
SWAPS	<input type="text"/>			
MOVIE	<input type="text"/>			

### For more help visit <https://www.deepfakes.club>

- Open this page again by clicking . Save or load settings by clicking and in the top left corner.

- VIDEO A** will select video clip A.
- IMAGES A** will collect frames from video clip A.
- FACES A** will extract and align faces from image set A.
- MODEL** will train a new or existing model.
- SWAPS** will convert image set A into face B by default.
- MOVIE** will combine the swapped images into a movie.

- Use to select the directory to store results from each step. You can also load a directory with pre-existing results for the next step.

- Inspect results by clicking to open the directory with Explorer.

- Open the options menu by clicking . You can also enter custom commands from there.

- Commands with empty paths will use the default directories with

# OpenFaceSwap





# Research question

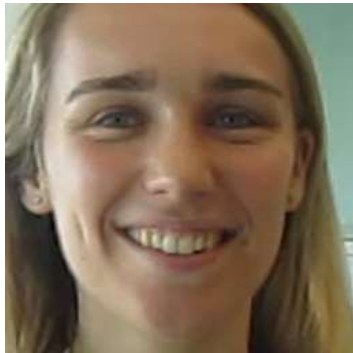
Can Deepfakes be distinguished from authentic videos with the use of:

1. Convolutional Neural Networks (CNN)
2. Discrete Cosine Transformation (DCT) coefficients analysis
3. Photo Response Non Uniformity (PRNU) analysis

- A pioneering study into the detection of Deepfakes.

# Dataset

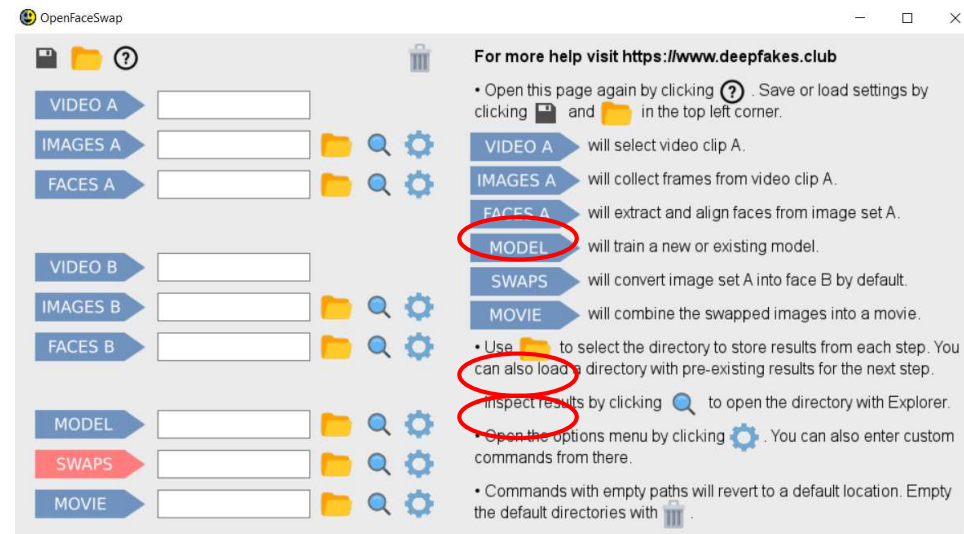
- 16 Deepfakes, 10 authentic videos
- Average video length = 29 seconds
- Frames extracted with FFmpeg
- Three actors used to create dataset





# Discrete Cosine Transformation (DCT): Method

- DCT are used in JPEG compression - can leave traces
  - Can be used to see whether the file has been saved more than once
- Authentic frame vs Swap frame
- Authentic frame vs Deepfake movie frame (extracted by FFmpeg)



# DCT: Results

Deepfake frame before video assembly

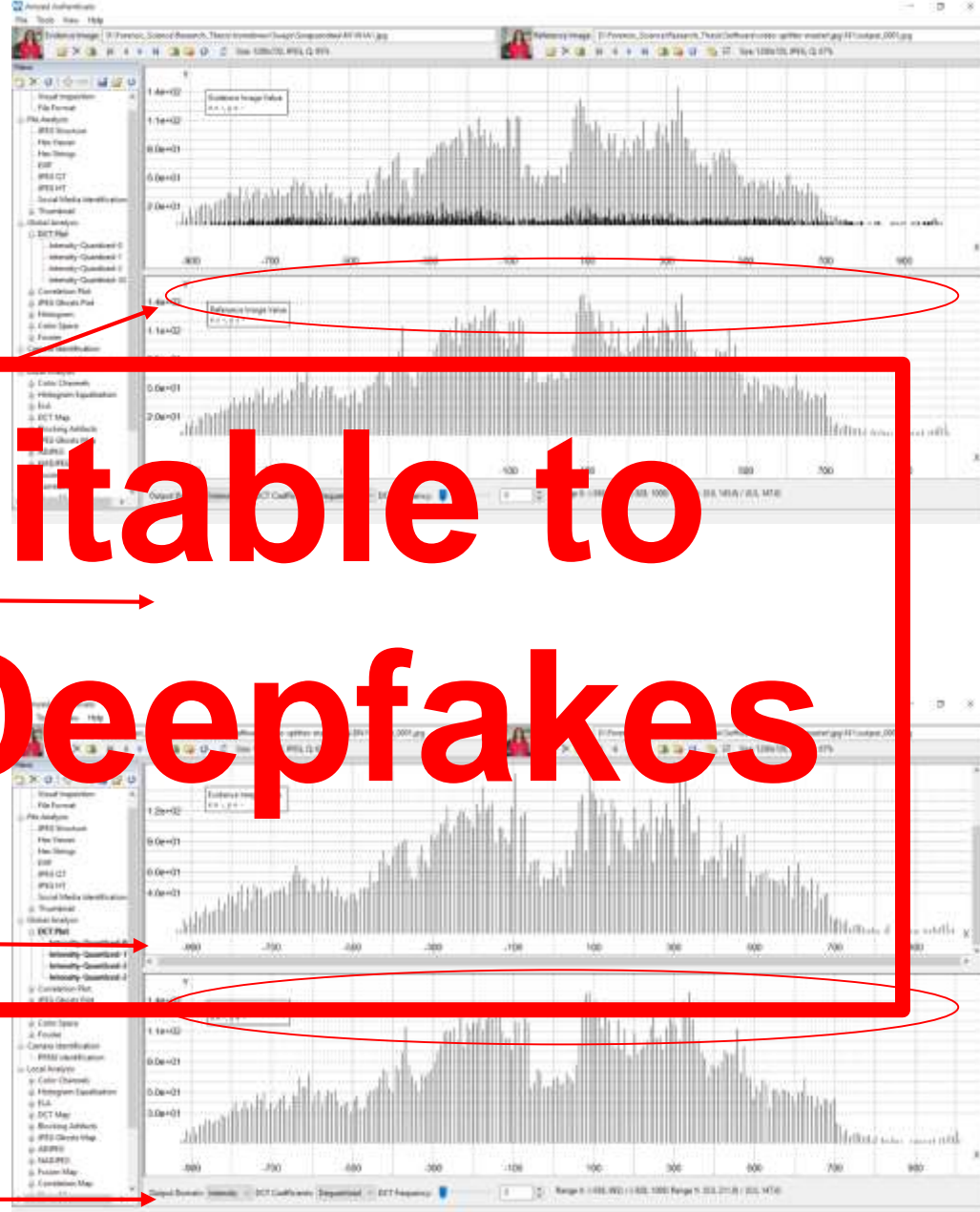
Authentic frame

Deepfake frame extracted from video after assembly

No sign of double compression

Authentic frame

**Not suitable to detect Deepfakes**





# Convolutional Neural Network (CNN): Method

- CNN are AI which take an input and learn to eg. classify it, without a human telling them how to do so.
- CNN based on GoogLeNet
- Classify frames from dataset as 'Natural' or 'Deepfake'
- 60/20/20
- Model trained for max. 100 epochs



# CNN: Results

- Classified all frames as 'Natural' (authentic), or all as 'Deepfake'
- Changing regularisation methods had no effect.

## GoogLeNet 5 Image Classification Model

### Summary

Top-1 accuracy  
62.13%

Top-5 accuracy  
100.0%

- Initialize
- Running
- Done a
- (Total: 5)

Infer Mode

Notes

None

### Confusion matrix

	Deepfake	Natural	Per-class accuracy
Deepfake	0	837	0.0%
Natural	0	1373	100.0%

### All classifications

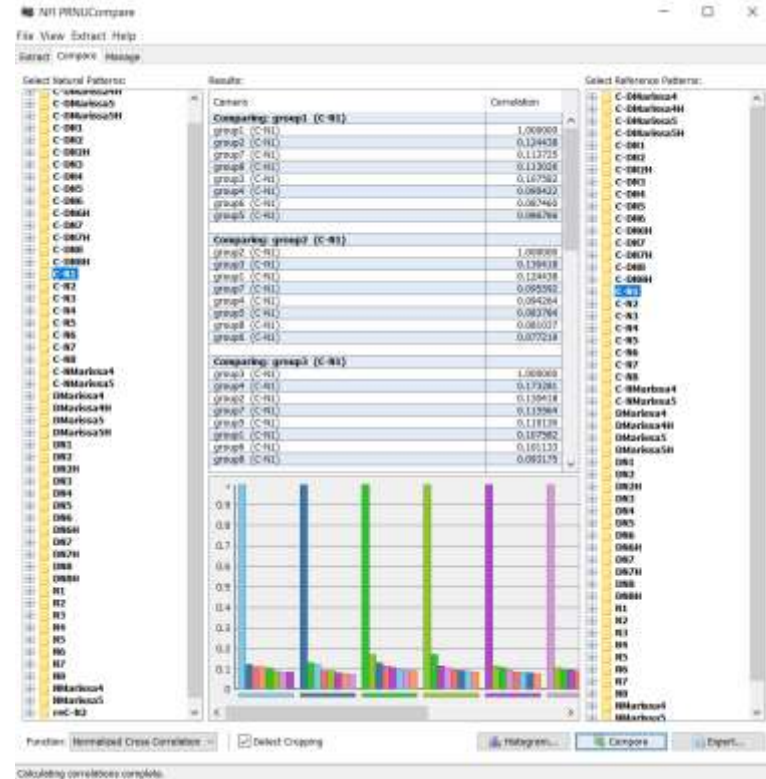
	Path	Ground truth	Top predictions
1	/test/Natural/\$E04480.png	Deepfake	Natural 69.47%
2	/test/Natural/\$E04175.png	Deepfake	Natural 67.64%

# CNN: Conclusions

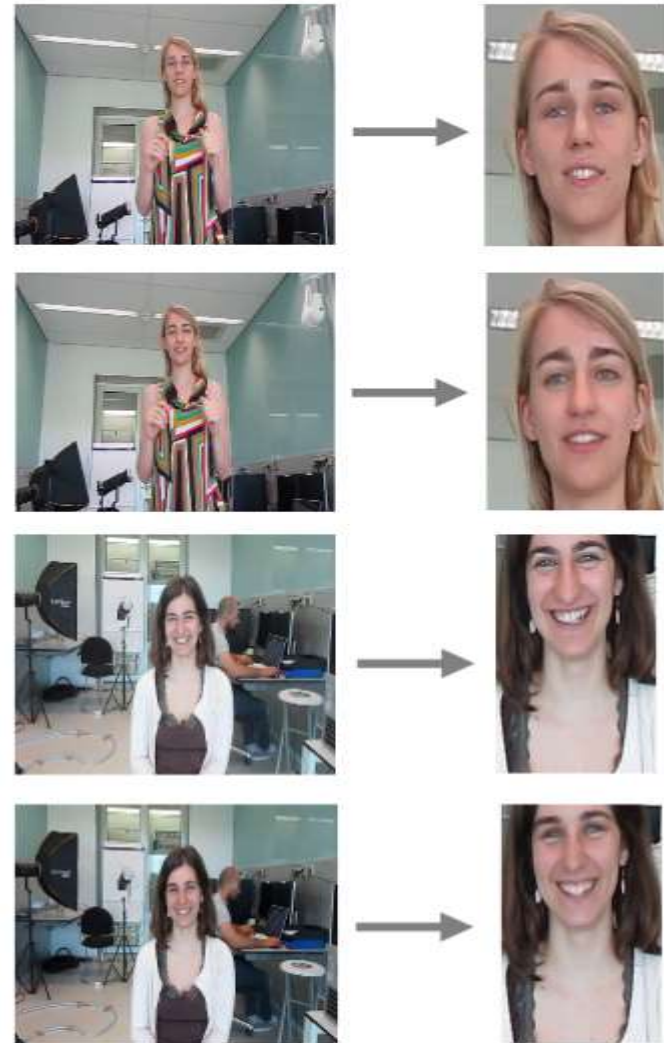
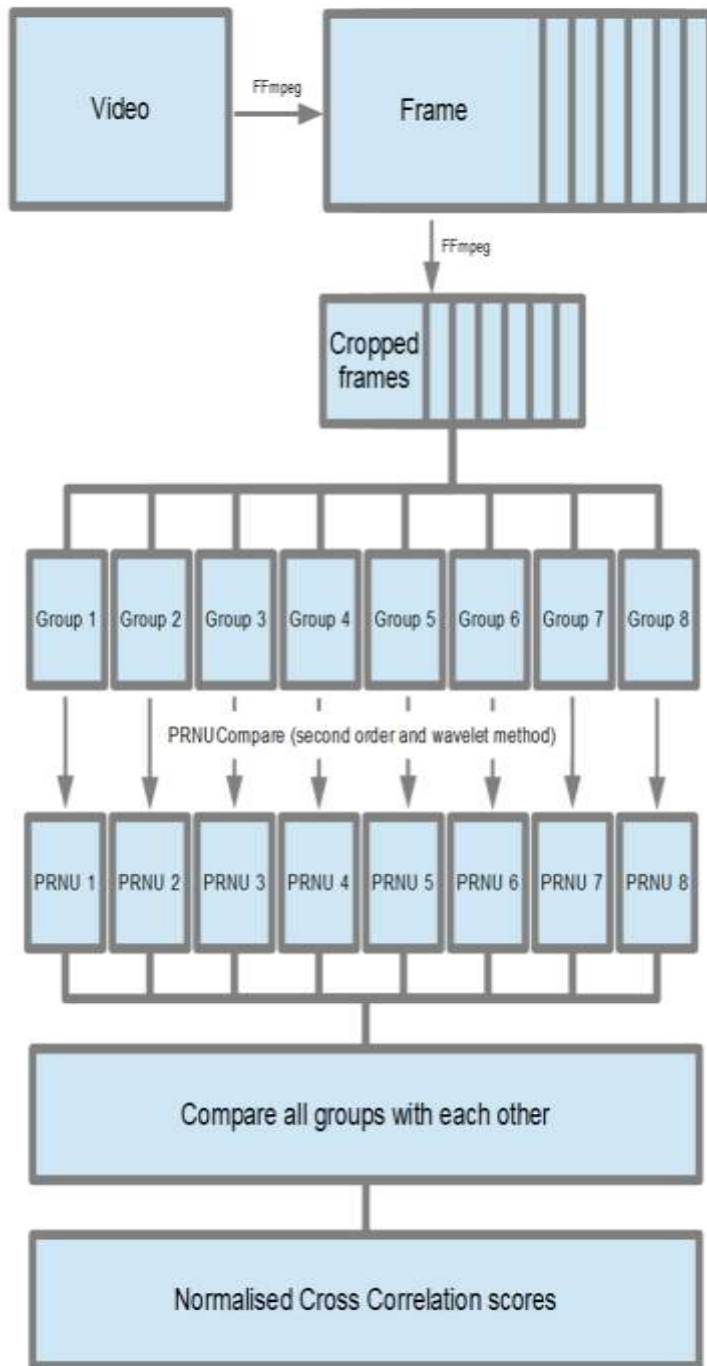
- Persistent overtraining
- Dataset with even number of Deepfake and Natural frames may improve learning
- Cannot conclude that CNNs are unsuitable
- General Adversarial Networks (GANs) may be more suitable.

# Photo Response Non Uniformity (PRNU) analysis: Method

- 'The fingerprint of the digital camera'
- Manipulation can alter PRNU pattern
- **Deepfake PRNU pattern less consistent throughout video compared to Authentic?**
- Second order and Wavelet method
  - Second wave = faster
  - Wavelet = more reliable
- Cropped and uncropped
  - Increase % variability of PRNU



# PRNU analysis: Method





# PRNU analysis: Conclusions

Experiment	P-value variance in normalised cross correlation scores	P-value mean normalised cross correlation scores
Second order cropped	0.593	$5.21 * 10^{-5}$
Second order uncropped	0.303	0.002
Wavelet cropped	0.041	0.188
Wavelet uncropped	0.852	$3.23 * 10^{-4}$

- Second order > Wavelet method
  - Takes less time
  - Stronger correlation
  - More reliable correlation
- Second order cropped > Second order uncropped
  - Stronger correlation

# Limitations of study

- Dataset
  - Small
  - Imbalanced
  - One camera
- CNN
  - Only CNN based on GoogLeNet
  - Imbalanced dataset
- PRNU analysis
  - Cropping method not suitable for videos with large movements
  - All results from one camera's PRNU pattern



**WARNING**

# Further Research

- PRNU
  - Confirm correlation
  - Likelihood ratios
  - Effect different cameras
  - Effect camera software (phone apps etc)
  - Is second order uncropped be sufficient
- CNN
  - General Adversarial Networks (GAN)
  - Balanced dataset



Thank you for your time

Questions [zeno@holmes.nl](mailto:zeno@holmes.nl) /  
[andrea@holmes.nl](mailto:andrea@holmes.nl)



## Discussion

Rafferty said: “Cost-cutting and outsourcing has put the administration of justice at risk ... I don’t think it’s bad faith by the police. They have been under-resourced. They are swamped. In some of my cases it’s the police who have revealed material that’s helpful to the defence.”

Collie, the head of Discovery Forensics in London who mainly works for defendants, said: “The odds are stacked against the defence in many ways. We rarely get access to the actual piece of equipment. In the past I could go to the police station and see a phone or a computer and physically check it’s the right piece. Now everything comes prepackaged and is handed over on a hard drive or USB stick.”



## Collapsed rape prosecutions

### December: Liam Allan

The first case to be abandoned due to the failure by police to hand over crucial digital evidence was that of London student Liam Allan, 22, in December. Allan was charged with 12 counts of rape and sexual assault, but his trial was abandoned after police were ordered to hand over phone records that should have already been provided to the defence.

### December: Isaac Itiary

Shortly before Christmas, an alleged child rapist, Isaac Itiary, 25, was cleared at Inner London crown court when the prosecution offered no evidence. Material recovered from the phone of the complainant by police was only handed over to defence lawyers shortly before it was due to come to trial.

### January: Oliver Mears

In January, Oliver Mears, 19, a student at Oxford University, was charged with the rape of a 16-year-old girl in 2015 following



# Challenges



- Explain Deep Learning in court
- Bias in Model
- Training of users

## Shielding

- Anti forensic software
- Encryption
- Darknets
- crypto currencies





# Questions



11