

Facial Recognition Performance and Its Measurement



PRESENTED BY:

Patrick Grother

National Institute of Standards and Technology

MODERATED BY:

Stephen Redifer

2020-09-24



HDIAAC

Homeland Defense & Security
Information Analysis Center

DISTRIBUTION STATEMENT A. Approved for public release; distribution is unlimited.

HDIAC is sponsored by the Defense Technical Information Center (DTIC). Any opinions, findings and conclusions or recommendations expressed in this material are those of the author(s) and do not necessarily reflect the views of the Defense Technical Information Center.

info@hdiac.org
<https://www.hdiac.org>



Face Recognition Performance and its Measurement

Patrick Grother

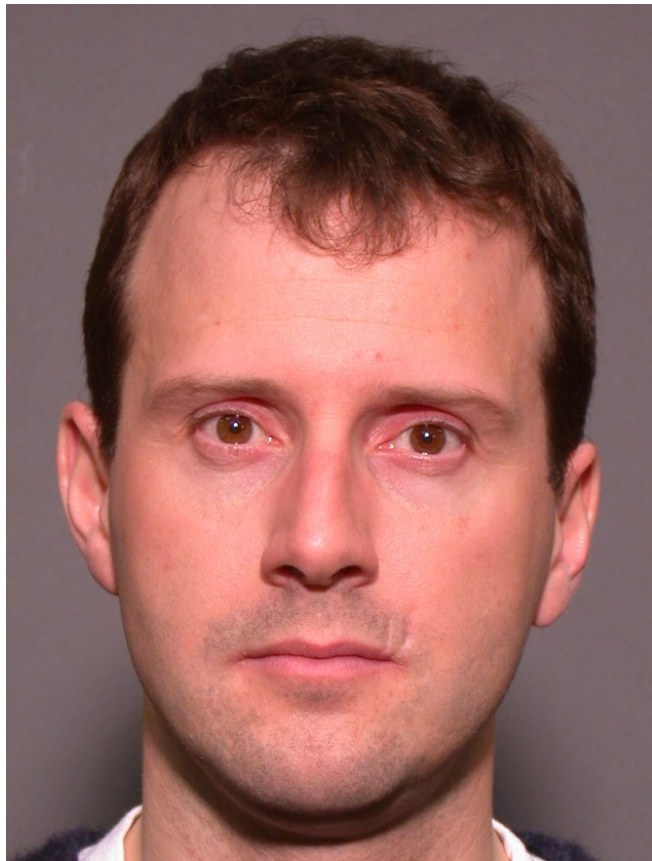
NIST

U. S. Department of Commerce

Homeland Defense and Security Information Analysis Center

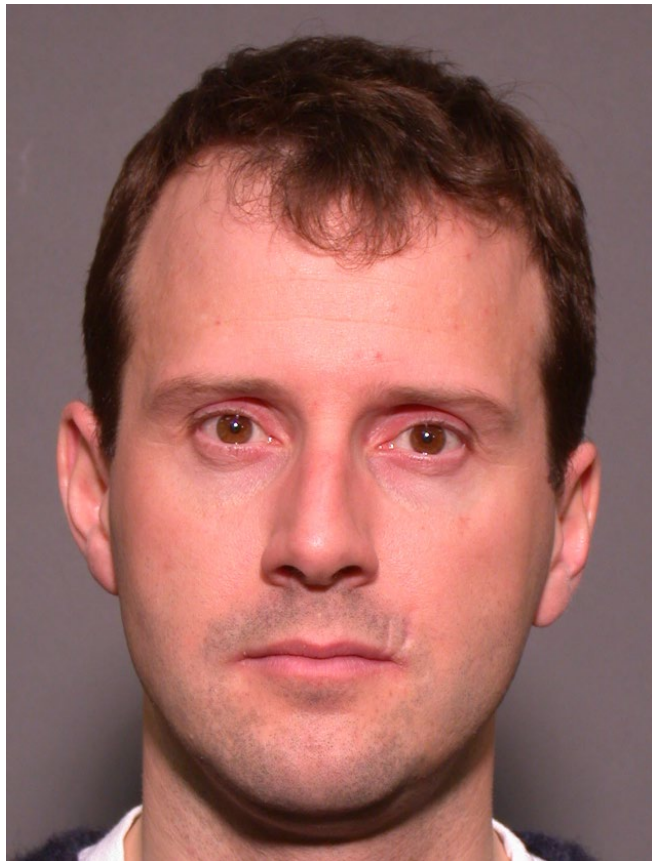
September 24, 2020

What's in a face?



How many biometrics here?

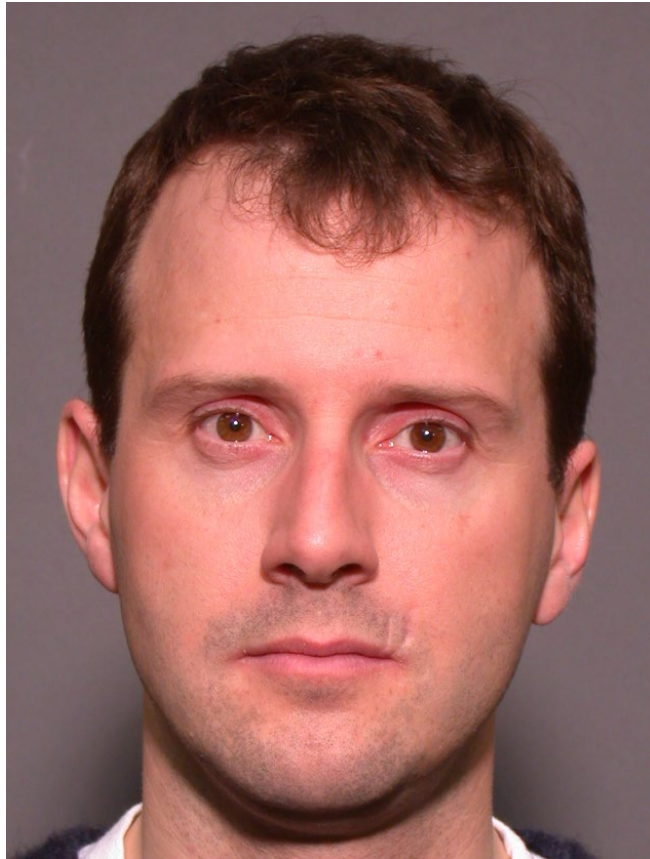
What's in a face?



How many biometrics here?

1 Face

What's in a face?

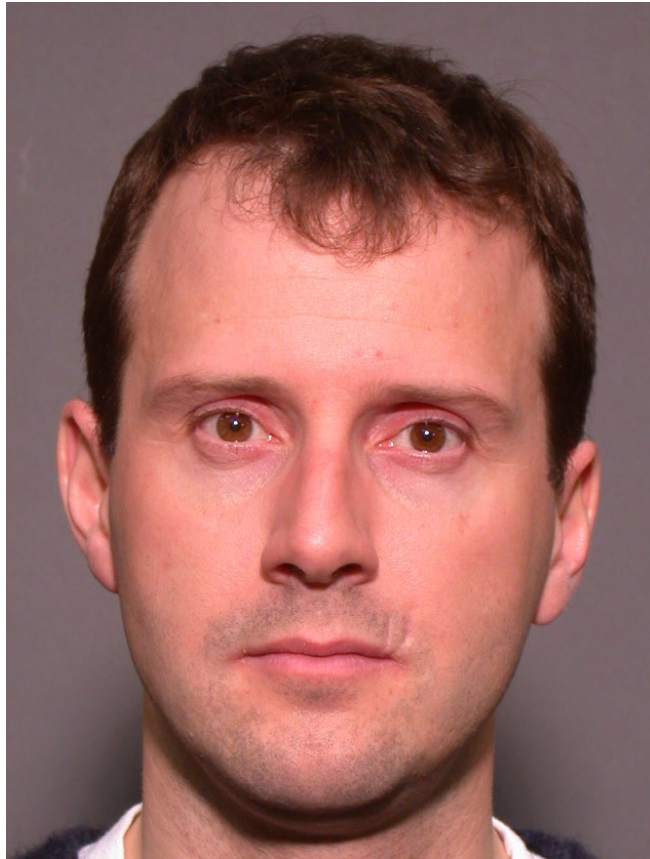


How many biometrics here?

- 1 Face
- 2 Irides + periocular



What's in a face?



How many biometrics here?

1 Face

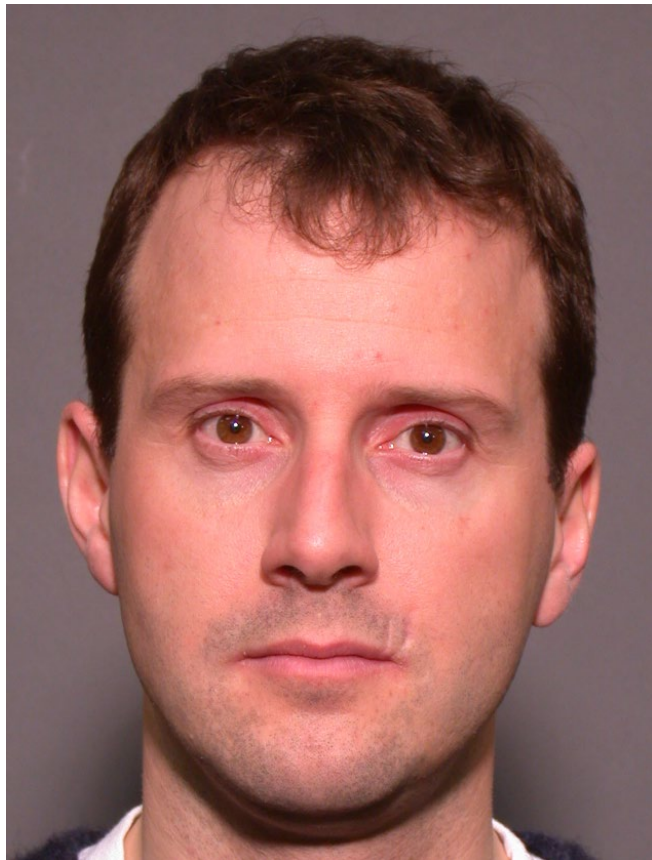
2 Irides + periocular

3 Skin texture

<https://patents.google.com/patent/US7369685B2/>



What's in a face?



How many biometrics here?

1 Face

2 Irides + periocular

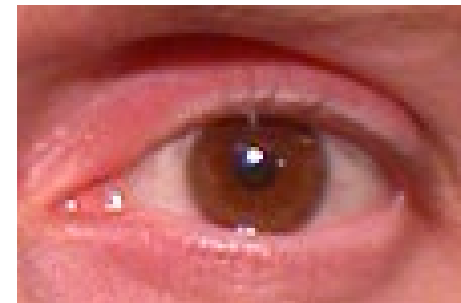
3 Skin texture | <https://patents.google.com/patent/US7369685B2/>

4 Head shape

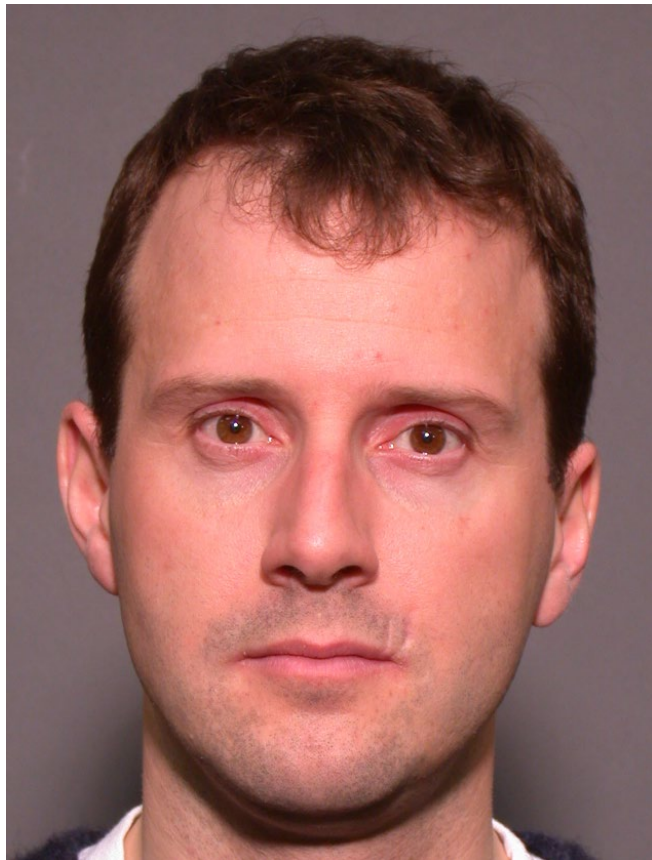
5 Ears

6 Scars

Human review: See ASTM E3149
*Standard Guide for Facial Image Comparison Feature
List for Morphological Analysis*

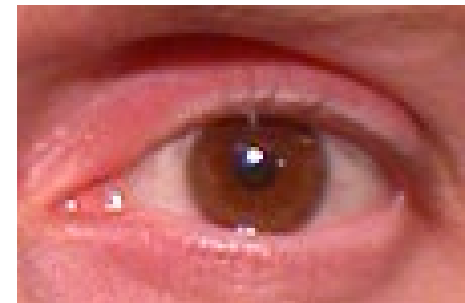


What's in a face?



How many biometrics here?

- 1 Face
- 2 Irides + periocular
- 3 Skin texture | <https://patents.google.com/patent/US7369685B2/>
- 4 Head shape | Human review: See ASTM E3149
Standard Guide for Facial Image Comparison Feature List for Morphological Analysis
- 5 Ears
- 6 Scars
- 7 **Anything else unique**
 - Short + long wave infrared
 - Hyperspectral
 - 3D



The Afghan Girl



STEVE MCCURRY



<https://www.nationalgeographic.com/magazine/2002/04/afghan-girl-revealed/>
c. National Geographic, photographic portrait by journalist Steve McCurry, 1984

Face authentication: Closed system



<https://spectrum.ieee.org/tech-talk/consumer-electronics/gadgets/new-samsung-galaxy-s8-unlocks-with-facial-recognition-iris-scanning>



<https://www.macrumors.com/2017/10/25/apple-reduced-face-id-accuracy-iphone-x/>



<https://support.apple.com/en-us/HT208109>

Face Recognition: How? By comparing faces



<https://securitytoday.com/articles/2018/02/27/us-border-patrol-unable-to-validate-epassport-data.aspx>

- Same identity?
- Different identity?



<https://www.thalesgroup.com/en/markets/digital-identity-and-security/government/eborder/eborder-abc>



Georgetown Law. Center on Privacy + Technology

<https://www.airportfacescans.com/>

Figure 2: A traveler has his face scanned as a Customs and Border Protection agent provides instruction. (Photo: Associated Press, all rights reserved)



https://en.wikipedia.org/wiki/FIPS_201



Face Recognition: How? By comparing faces



<https://securitytoday.com/articles/2018/02/27/us-border-patrol-unable-to-validate-epassport-data.aspx>

- Same identity?
- Different identity?



<https://www.thalesgroup.com/en/markets/digital-identity-and-security/government/eborder/eborder-abc>



Georgetown Law. Center on Privacy + Technology

<https://www.airportfacescans.com/>

Figure 2: A traveler has his face scanned as a Customs and Border Protection agent provides instruction. (Photo: Associated Press, all rights reserved)



https://en.wikipedia.org/wiki/FIPS_201

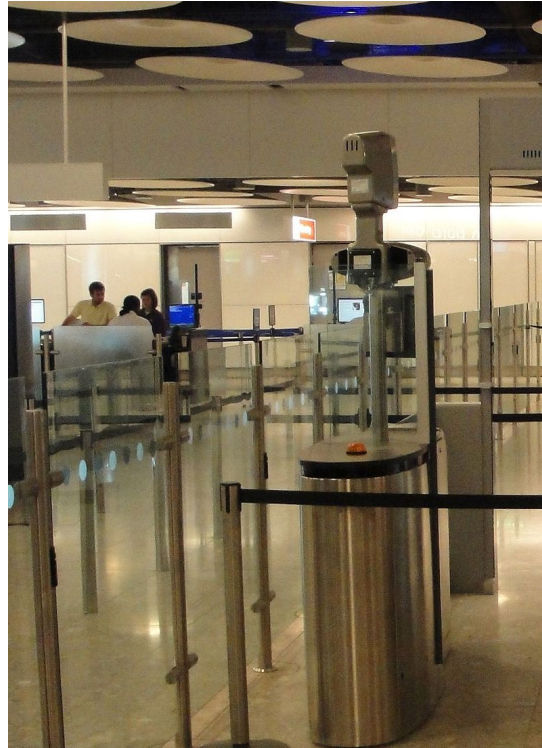
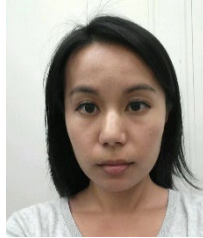


Source: FRVT staff and sister, with permission

Inbound border crossing using passport verification

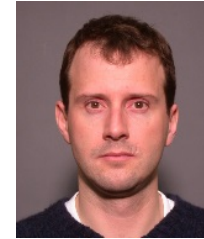


LIVE

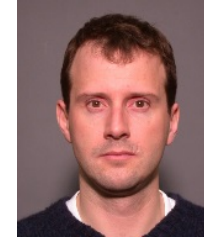


https://en.wikipedia.org/wiki/EPassport_gates
CC BY 2.0. File:Heathrow Terminal 5 ePassport gates.jpg
Created: 16 July 2010

CHIP



Wait for luggage



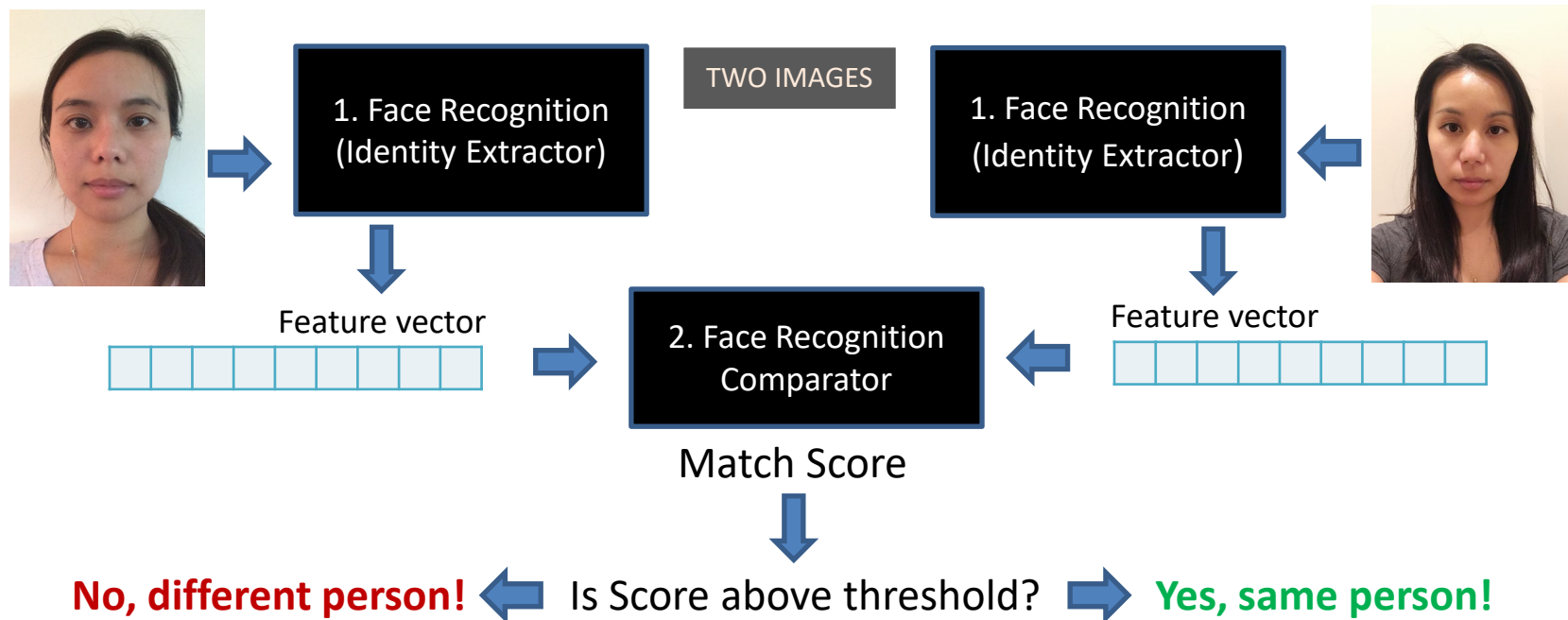
See an officer!

Two factor authentication:

1. Something you have:
2. Something you are:

Possession of passport
Successful recognition of a biometric

Face recognition: How?

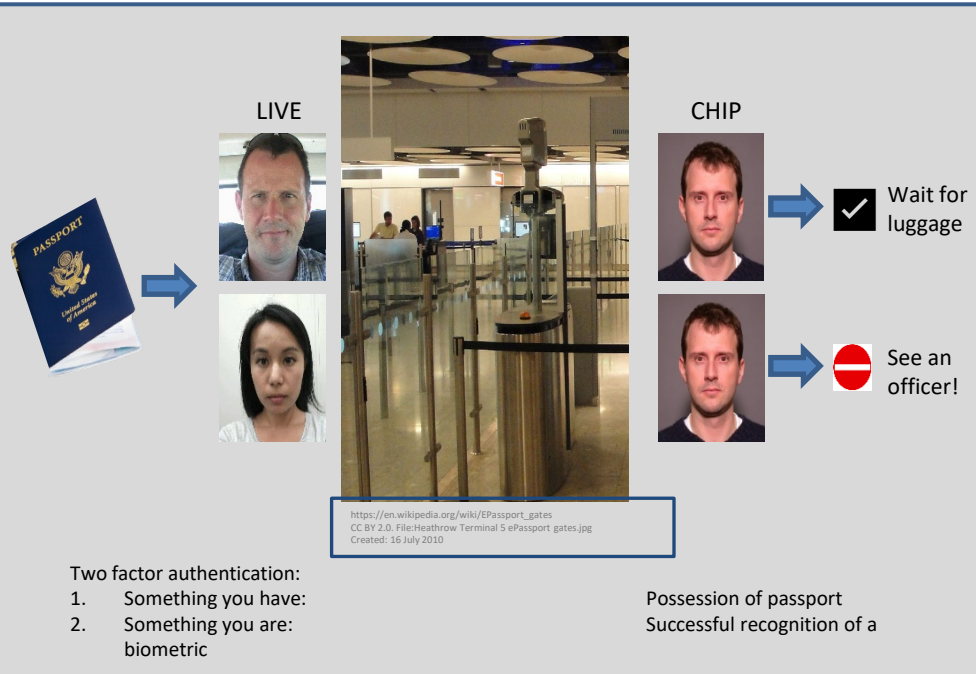


- DCNNs
- ML / AI
- Not commoditized
- Trade secrets

- Templates aka feature vectors
- 0.2 - 4KB, 2KB is most common
 - 0.1 to 1 second on CPU

- Templates
- Templates are reversible
 - Images retained

FR in operations: Passport verification at a border



1. No central database
2. Two images involved: live capture and chip image
3. Trusted passport?
 - Digital signature
 - Morphed image
4. Error and consequences
 - False Accept → Border security
 - False Negative → Inconvenience

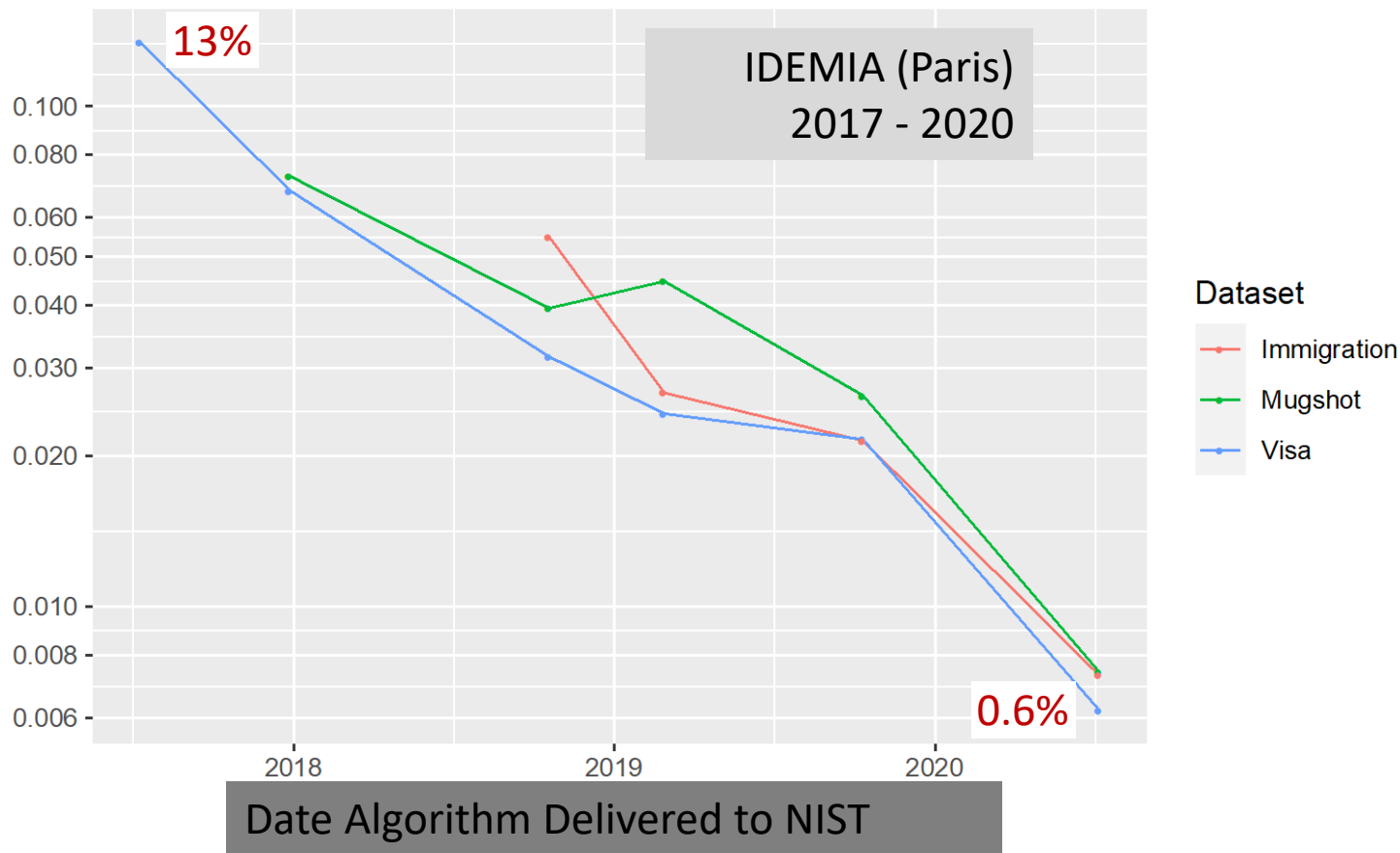
FRVT 1:1 Leaderboard 2020-07-27

Developer	VISA Photos FNMR @ FMR ≤ 0.000001	MUGSHOT Photos FNMR @ FMR ≤ 0.00001	MUGSHOT Photos FNMR @ FMR ≤ 0.00001 DT ≥ 12 YRS	VISABORDER Photos FNMR@ FMR ≤ 0.000001	BORDER Photos FNMR @ FMR = 0.000001	WILD Photos FNMR@ FMR ≤ 0.00001
sensetime-003	0.0027 ⁽³⁾	0.0027 ⁽¹⁾	0.0027 ⁽¹⁾	0.0051 ⁽⁶⁾	0.0100 ⁽⁷⁾	0.0355 ⁽⁴⁵⁾
deepglint-002	0.0027 ⁽²⁾	0.0032 ⁽⁷⁾	0.0033 ⁽²⁾	0.0043 ⁽²⁾	0.0084 ⁽³⁾	0.0301 ⁽¹⁾
paravision-004	0.0046 ⁽⁷⁾	0.0030 ⁽⁴⁾	0.0036 ⁽³⁾	0.0091 ⁽¹⁸⁾	0.0188 ⁽²⁷⁾	0.0311 ⁽¹⁶⁾
visionlabs-008	0.0036 ⁽⁴⁾	0.0031 ⁽⁶⁾	0.0040 ⁽⁴⁾	0.0045 ⁽³⁾	0.0079 ⁽¹⁾	0.0308 ⁽¹⁰⁾
...						
toshiba-003	0.0214 ⁽⁶⁴⁾	0.0085 ⁽⁴¹⁾	0.0131 ⁽⁴⁰⁾	-	0.0241 ⁽³⁷⁾	0.0321 ⁽²⁶⁾
fujitsulab-000	0.0212 ⁽⁶³⁾	0.0091 ⁽⁴⁵⁾	0.0133 ⁽⁴¹⁾	0.0251 ⁽⁷¹⁾	0.4200 ⁽¹⁰⁵⁾	0.0481 ⁽⁷³⁾
asusaics-000	0.0209 ⁽⁶²⁾	0.0085 ⁽³⁹⁾	0.0134 ⁽⁴²⁾	0.0143 ⁽³⁸⁾	0.7189 ⁽¹¹²⁾	0.0332 ⁽³⁵⁾
cogent-004	0.0116 ⁽³³⁾	0.0096 ⁽⁴⁹⁾	0.0134 ⁽⁴³⁾	0.0157 ⁽⁴¹⁾	0.0325 ⁽⁵⁴⁾	0.0436 ⁽⁶⁶⁾

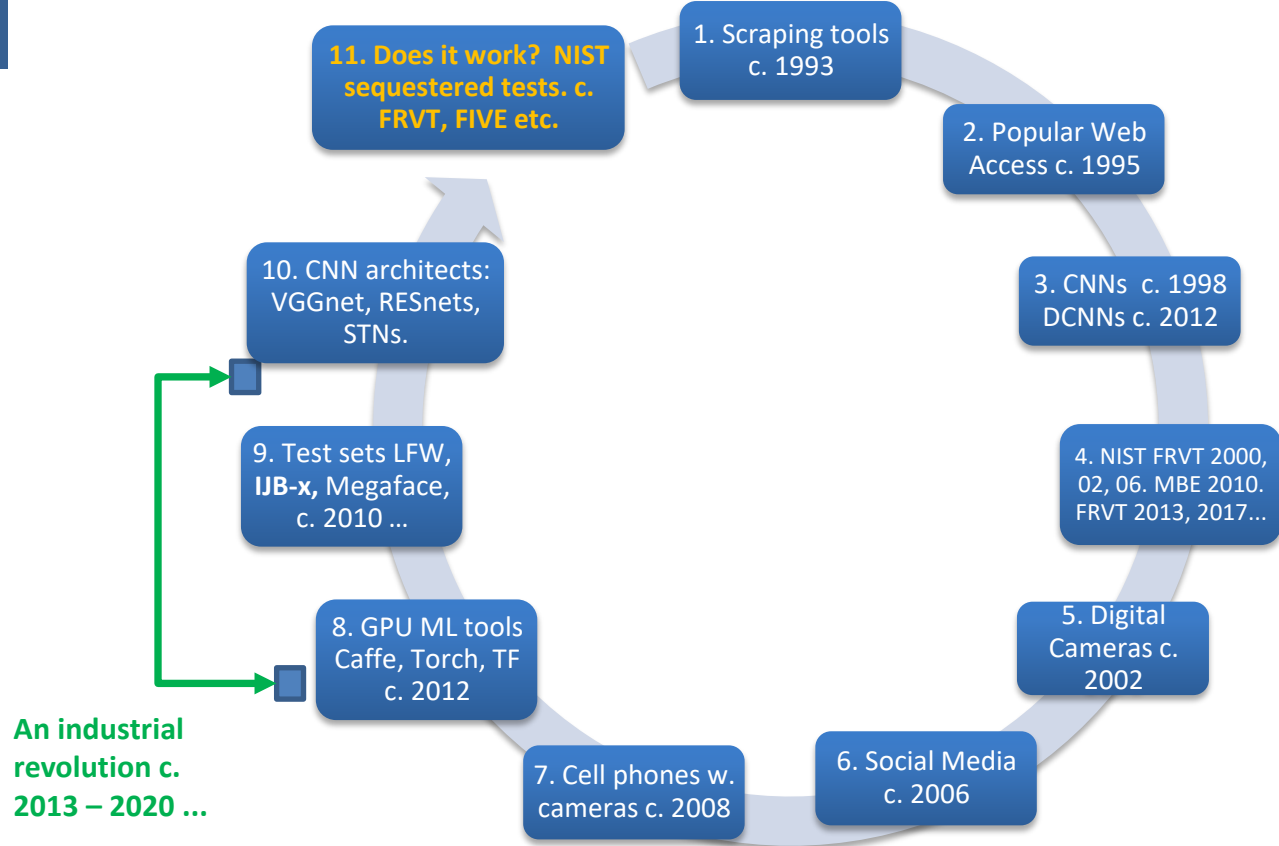
Accuracy Gains: Typical Example

False
Rejection

FNMR at
 $FMR = 10^{-6}$



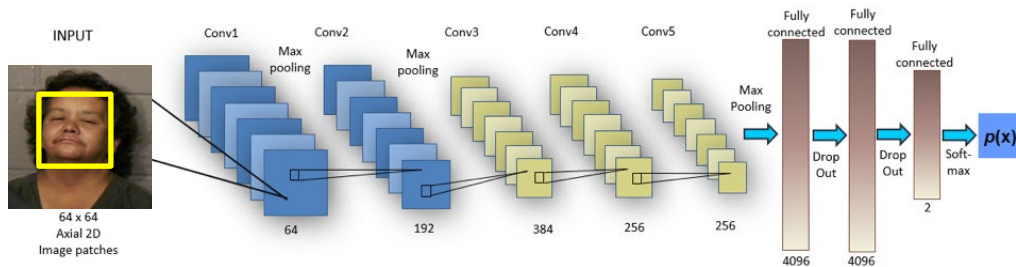
Enablers of Better Face Recognition



Black box: What is a DCNN?

$$F(\mathbf{x}) = F_N(F_{N-1}(\dots F_2(F_1(\mathbf{x}, \mathbf{w}_1), \mathbf{w}_2) \dots, \mathbf{w}_{N-1}), \mathbf{w}_N)$$

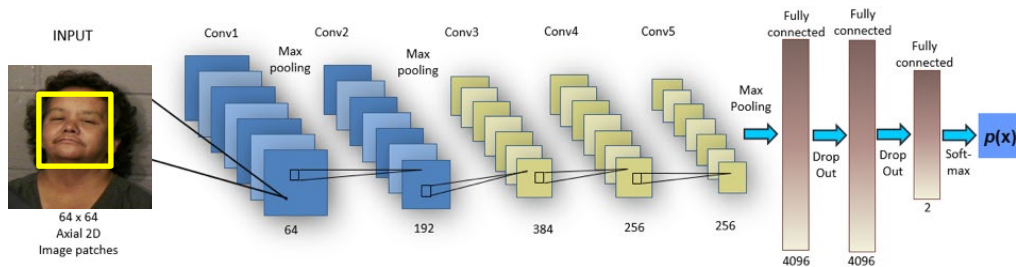
- » CNN is a composed function, F , implementing local (image) filters
- » Operating on an image, \mathbf{x}_1 , input to the first layer
 - Dimensions are $W \times H \times K$
- » Producing intermediate feature maps, \mathbf{x}_n , $1 < n \leq N$
- » Each layer has a function, F_n , which perform various operations and are handcrafted
- » Each layer has parameters, \mathbf{w}_n , which are **learned from some training data**



Black box: What is a DCNN?

$$F(\mathbf{x}) = F_N(F_{N-1}(\dots F_2(F_1(\mathbf{x}, \mathbf{w}_1), \mathbf{w}_2) \dots, \mathbf{w}_{N-1}), \mathbf{w}_N)$$

- » CNN is a composed function, F , implementing local (image) filters
- » Operating on an image, \mathbf{x}_1 , input to the first layer
 - Dimensions are $W \times H \times K$
- » Producing intermediate feature maps, \mathbf{x}_n , $1 < n \leq N$
- » Each layer has a function, F_n , which perform various operations and are handcrafted
- » Each layer has parameters, \mathbf{w}_n , which are **learned from some training data**



Multimodal



Multisensor



Multi-instance
(contemporaneous)



Repeated-instance
(longitudinal)



Multiple algorithm

Score = Fusion [Algorithm_B(X,Y), Algorithm_A(X,Y)]

1:1 Authentication: Live-to-document



<https://www.thalesgroup.com/en/markets/digital-identity-and-security/government/eborder/eborder-abc>



https://en.wikipedia.org/wiki/FIPS_201



Georgetown Law. Center on Privacy + Technology
<https://www.airportfacescans.com/>

Figure 2: A traveler has his face scanned as a Customs and Border Protection agent provides instruction. (Photo: Associated Press, all rights reserved)

1:N Identification

Scalability to Large Populations

High volume applications

- Duplicate detection (passports, visa fraud, National ID)
- Casino persons of interest
- Aircraft boarding
- Surveillance

Human review usually infrequent

Low volume applications, with human review:

- Criminal investigation
- Clustering media

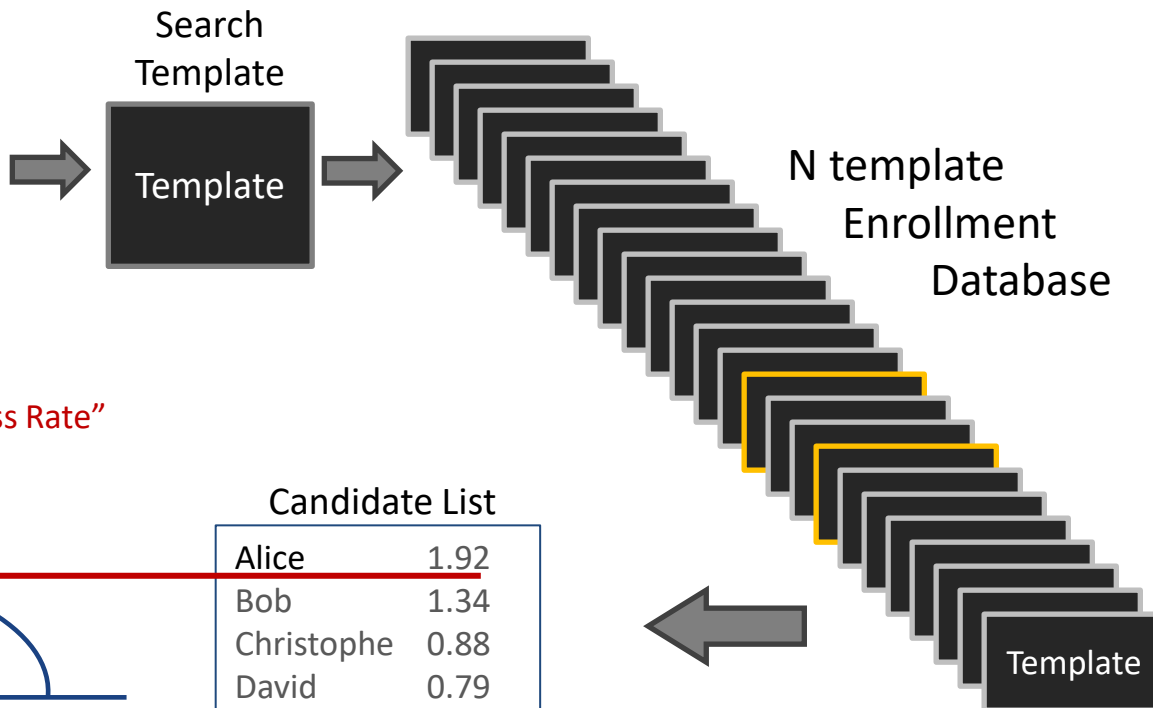
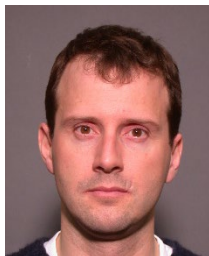


THIS IS **NOT** HOW FR WORKS. INSTEAD:

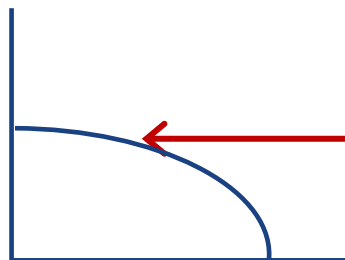
- An FR engine only knows people who are enrolled into it
- FR implements comparisons of new photos

1:N Search = N 1:1 comparisons (sometimes)

Biometric sample



FNIR, aka "Miss Rate"



FPIR

Aka False Alarm Rate

Candidate List

Alice	1.92
Bob	1.34
Christophe	0.88
David	0.79
Ernie	0.76

A demonstration of 1:N face recognition

» Enroll border crossing images

- 104.1 million
- 32.6 million people

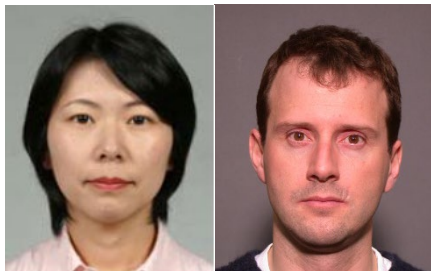


» Mated searches

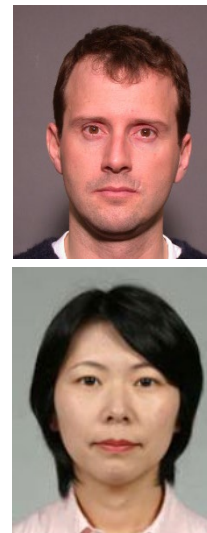
- 2.3 million “visa” APPLICATION images  FNIR, aka “miss rate”

» Non-mated searches

- 1.8 million “visa” APPLICATION images  FPIR, aka “false alarm rate”



A demonstration of 1:N face recognition



Step 1:

- Enrol N = 104 million photos, of 32.6 million people
- Images are examples, from NIST Special Database 32, representative of pose, illumination, compression

Step 2:

- Search with almost ISO compliant “visa” portraits

104 Million: “visa” to “border crossing” search accuracy

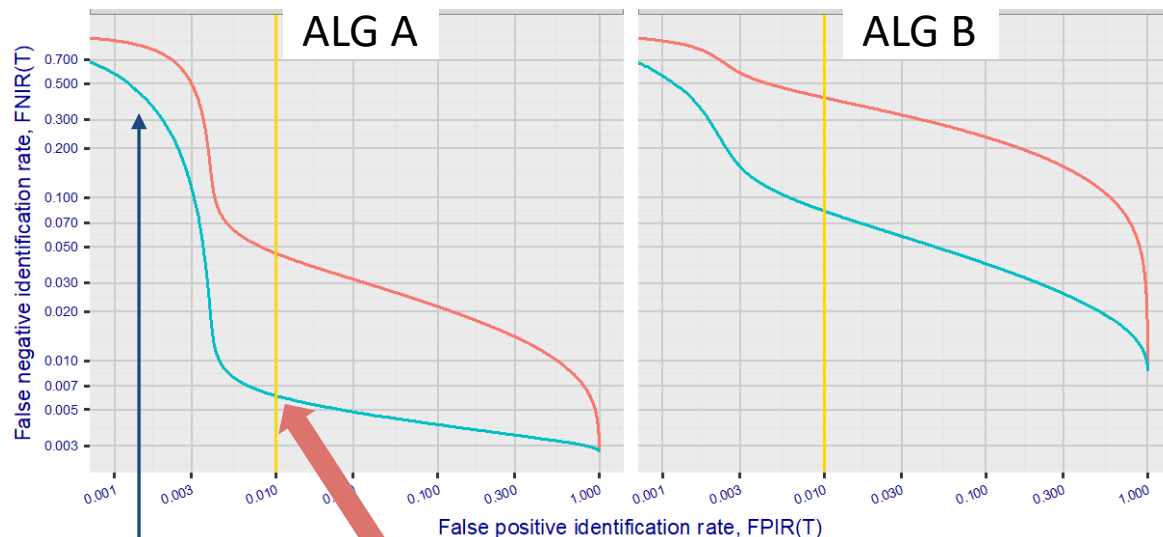
INVESTIGATION FALSE POS ID RATE = 100%	NEC-3 (2018-11) (0.7 + 1.1 seconds)	RankOne-006 (2019-06) (0.1 + 18 seconds)
Searches not returning ANY image of the correct person at rank 1	0.4%	2%
Searches not returning ALL images of the correct person in the top ranks	1.6%	11%

HIGH VOLUME, HIGH THRESHOLD IDENTIFICATION, FALSE POS ID RATE = 1%	NEC	Rank One
Searches not returning ANY image of the correct person above threshold	0.6%	8.3%
Searches not returning ALL images of the correct person above threshold	4.5%	41.0%

1:N search accuracy

Enroll N = 104 million ENTRY images; Search CIS Portraits

1. But version control matters:
 - ⇒ NIST eval vs. Productized
2. **Investigative search with N > 100M is possible, defensible**
3. Low FPIR is not attainable, limited by
 - ⇒ Unconsolidated IDs
 - ⇒ So do presence of twins > siblings > families



Miss rate: 0.6% ⇒ Hit rate: 99.4%

With threshold set so that only 1 in 100 non-mate search produces a false positive

FRVT 1:N Leaderboard 2020-08-12

Algorithm	Mugshot Mugshot N = 12000000	FBI FBI	Mugshot Mugshot N = 1600000	FBI FBI	Mugshot Webcam N = 1600000	FBI CBP	Mugshot Profile N = 1600000	FBI FBI	Visa Border N = 1600000	VISA AIRPORT
deepglint_001	-		0.0025 ⁽⁴⁾		0.0116 ⁽²⁾		0.7914 ⁽²⁴⁾		0.0051 ⁽¹⁾	
sensetime_003	0.0024 ⁽¹⁾		0.0015 ⁽¹⁾		0.0105 ⁽¹⁾		0.1953 ⁽⁴⁾		0.0067 ⁽²⁾	
nec_3	0.0031 ⁽²⁾		0.0021 ⁽³⁾		0.0149 ⁽³⁾		0.5136 ⁽¹³⁾		0.0070 ⁽³⁾	
paravision_005	0.0065 ⁽⁴⁾		0.0030 ⁽⁵⁾		0.0199 ⁽⁵⁾		0.2335 ⁽⁶⁾		0.0098 ⁽⁴⁾	
pixelall_004	0.0230 ⁽¹⁴⁾		0.0109 ⁽¹³⁾		0.0497 ⁽¹⁷⁾		0.9992 ⁽¹³⁶⁾		0.0227 ⁽⁵⁾	
microsoft_6	0.0184 ⁽⁹⁾		0.0086 ⁽¹⁰⁾		0.0298 ⁽⁹⁾		0.1174 ⁽¹⁾		0.0234 ⁽⁶⁾	
ntechlab_008	0.0218 ⁽¹¹⁾		0.0099 ⁽¹²⁾		0.0364 ⁽¹¹⁾		0.1998 ⁽⁵⁾		0.0284 ⁽⁷⁾	
idemia_007	0.0242 ⁽¹⁵⁾		0.0123 ⁽¹⁷⁾		0.0419 ⁽¹⁶⁾		1.0000 ⁽¹⁶⁸⁾		0.0350 ⁽⁸⁾	
rankone_009	0.0258 ⁽¹⁸⁾		0.0124 ⁽¹⁸⁾		0.0597 ⁽²⁴⁾		0.8180 ⁽²⁵⁾		0.0427 ⁽⁹⁾	
dermalog_007	0.1097 ⁽⁸¹⁾		0.0594 ⁽⁹⁵⁾		0.1202 ⁽⁸⁴⁾		0.9341 ⁽³⁹⁾		0.1027 ⁽¹⁰⁾	
gorilla_004	0.1109 ⁽⁸²⁾		0.0645 ⁽¹⁰⁷⁾		0.1317 ⁽⁹⁷⁾		0.8521 ⁽²⁶⁾		0.1059 ⁽¹¹⁾	

• Values are threshold-based FNIR at FPIR = 0.003

• <https://pages.nist.gov/frvt/html/frvt1N.html>

Performance

- » Massive expansion of industry
 - International markets + adoption
- » Massive gains in accuracy
 - Very accurate on high quality images
 - Better tolerance of poor image quality
 - Better tolerance of ageing (time lapse < 20 years)
 - Operate with larger databases
- » Accuracy varies greatly across the industry
 - China – EU – Japan – Russia – US
 - Buyer beware!
- » Some high volume applications (e.g. duplicate detection) require a high threshold for low false positives
 - Leads to higher false negatives
 - Image quality remains critical
- » Face-aware cameras
 - ISO/IEC 24358 camera capabilities

Limitations

- » Demographic differentials "bias"
 - False positive >> False negative
 - False negatives from poor quality photos
 - Large false positive variations by race
 - Higher false positives among women, elderly, young
 - Algorithm matters
 - Better accuracy → smaller inequities
 - Only some Chinese algorithms give false positive rates on Chinese faces similar to those in Caucasian
 - Some one-to-many algorithms mitigate differentials
 - "Know-your-algorithm"
- » Twins not separable (false positives)
- » Attacks
 - Easy to "steal" a face for impersonation
 - Systems may be deployed without attack detection
 - Morphing
 - Adversarial
- » Human review capability is poor

AGEING



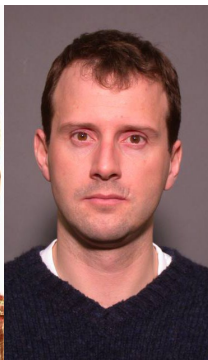
Images from presenter

Ageing

2002-08



2004-10



2010-05



2012



2013-08



2018-06

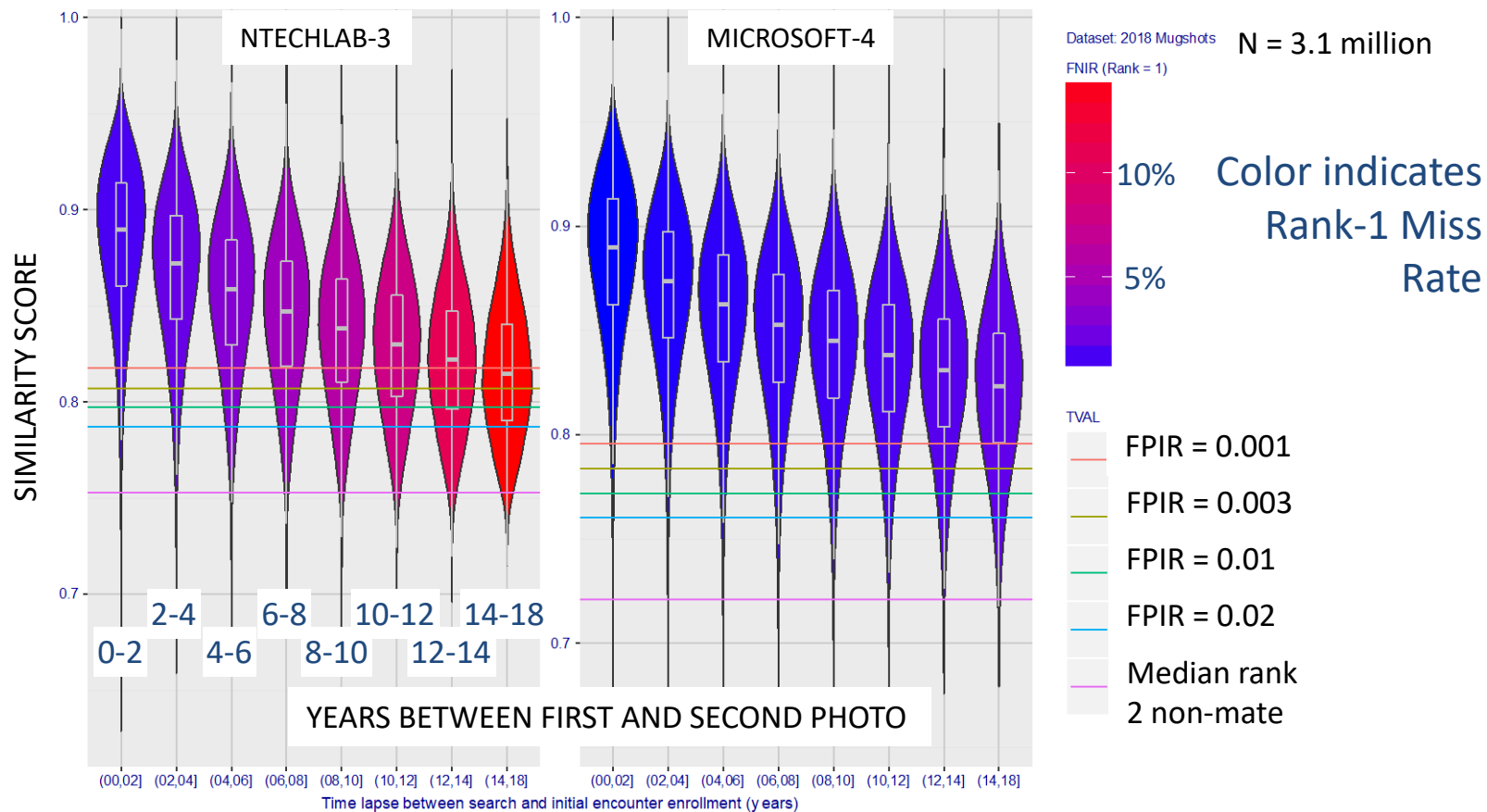


Images from presenter

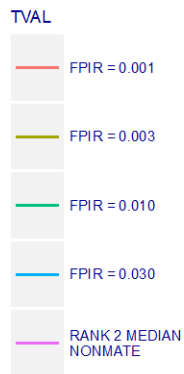
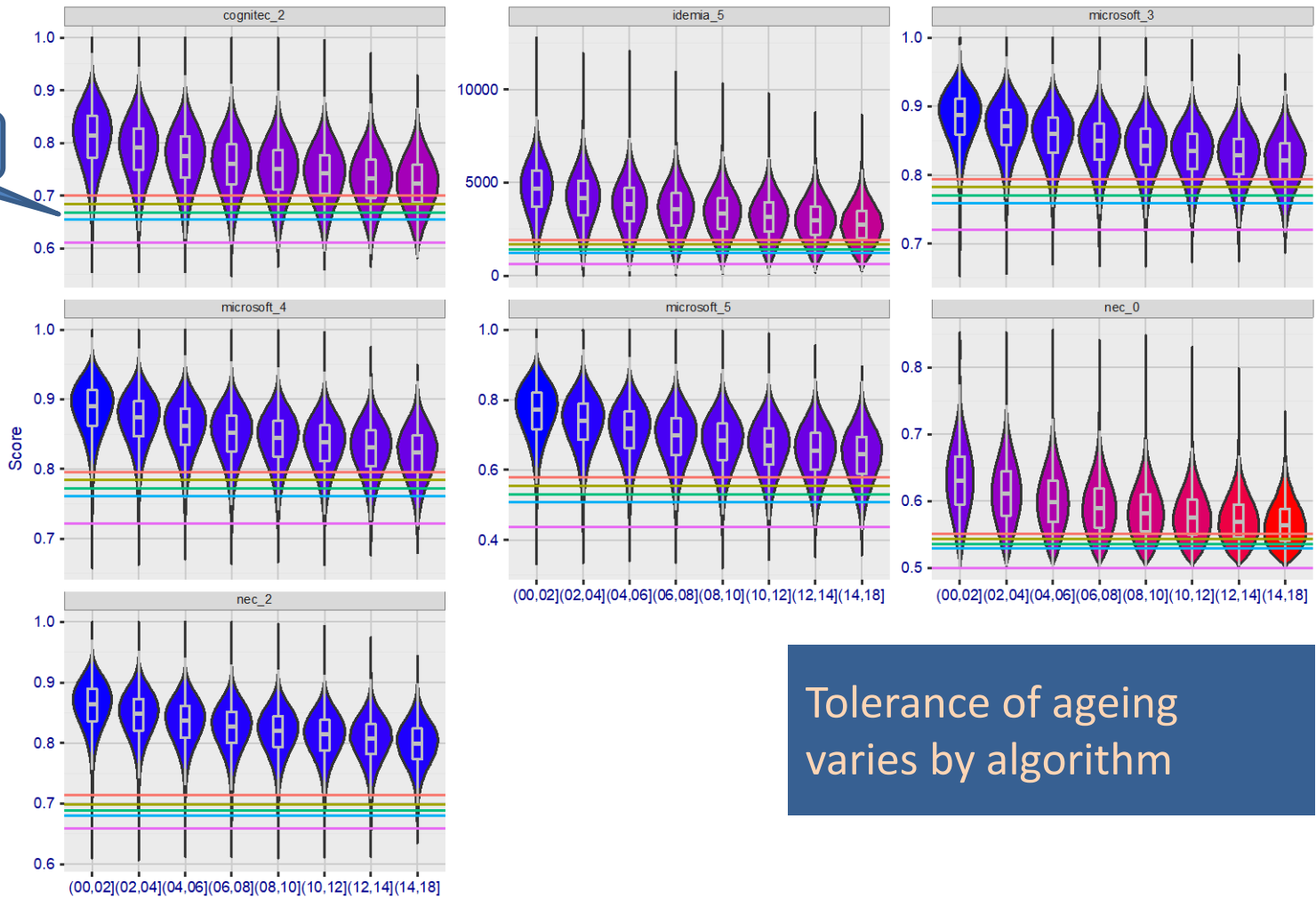


<https://www.bellingcat.com/news/uk-and-europe/2018/09/26/skrpal-suspect-boshirov-identified-gru-colonel-anatoliy-chepiga/>

Mate score distributions under ageing



Thresholds

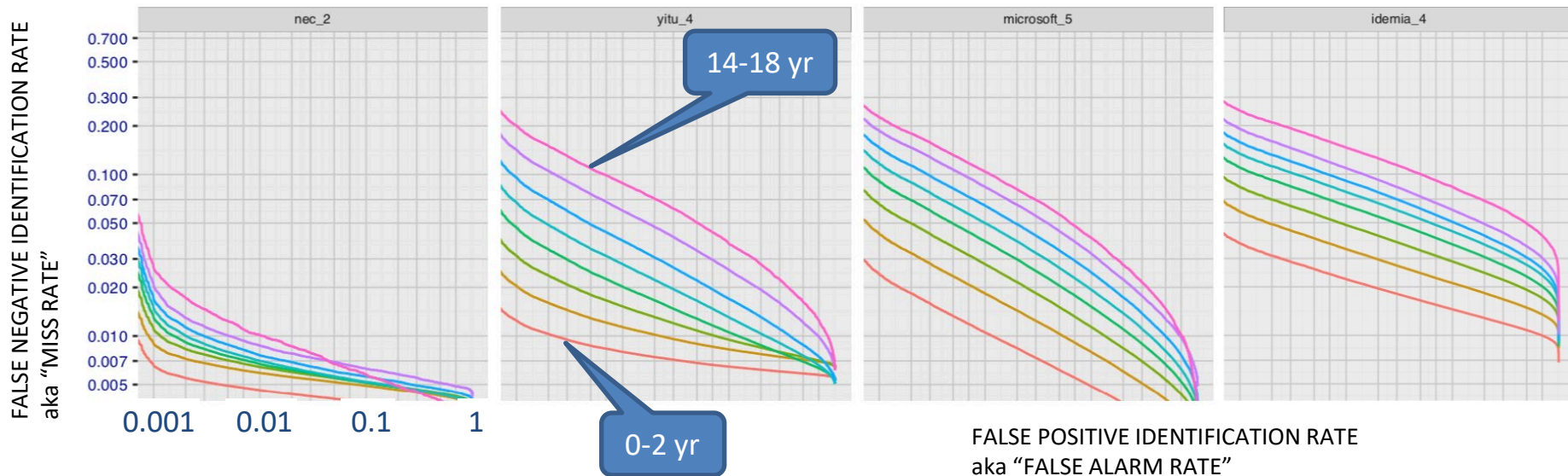


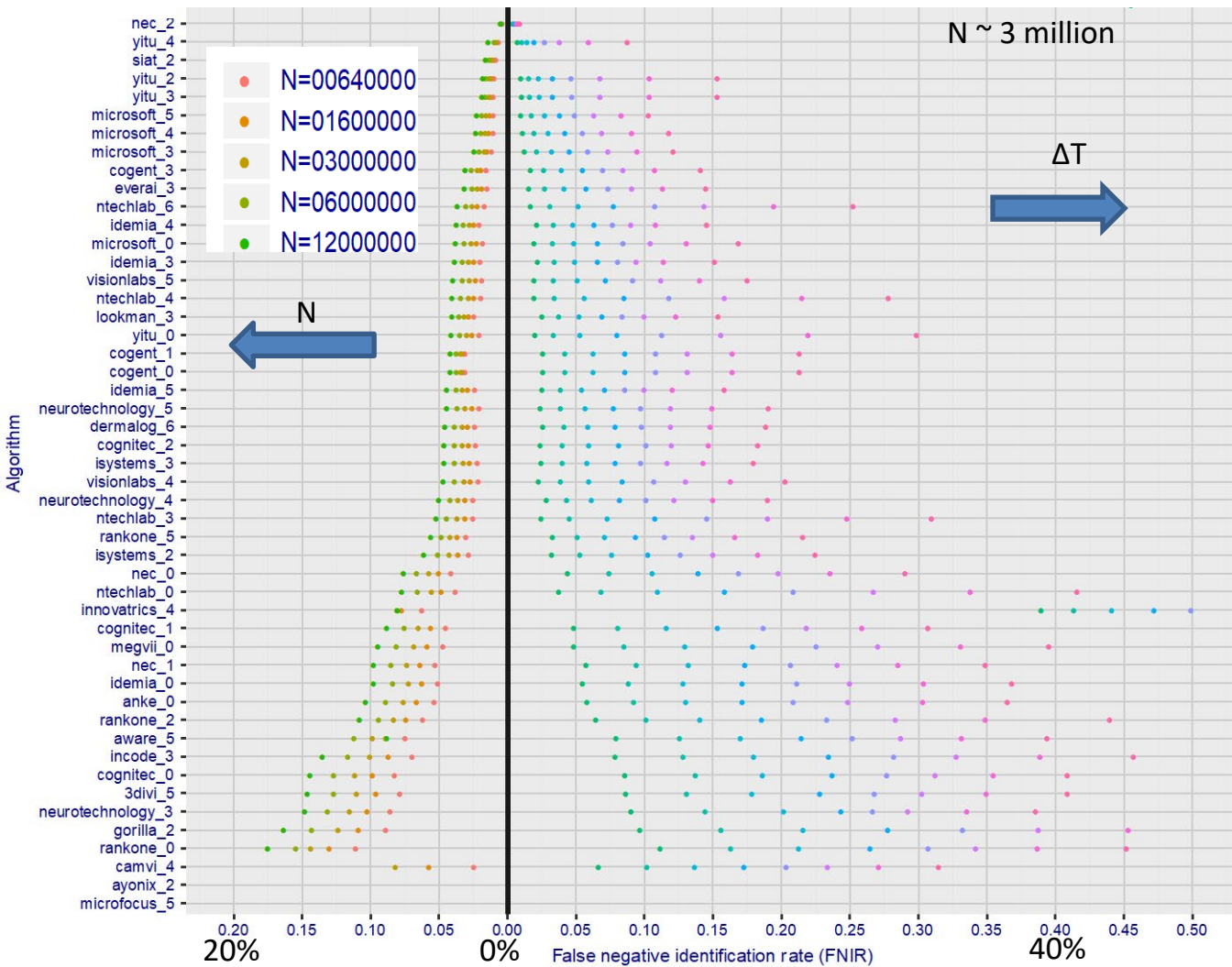
Dataset: 2018 Mugshots

Tolerance of ageing varies by algorithm

Time lapse between search and initial encounter enrollment (y ears)

Ageing: N = 3.1 million





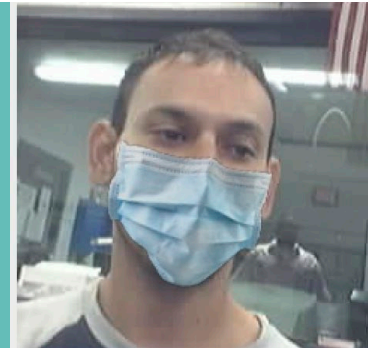
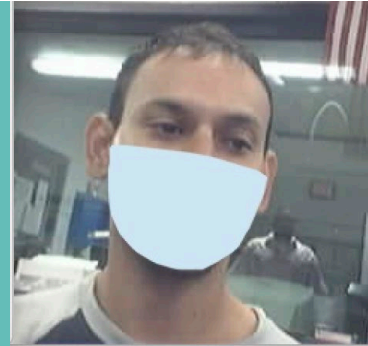
Performance in perspective: What matters more?

1. Algorithm
2. Population size
3. Ageing

- Years Lapsed (00,02]
- Years Lapsed (02,04]
- Years Lapsed (04,06]
- Years Lapsed (06,08]
- Years Lapsed (08,10]
- Years Lapsed (10,12]
- Years Lapsed (12,14]
- Years Lapsed (14,18]

Masks

What happens when you hide 40-70% of the face?



Synthetic masks

» NIST will vary

- Shape, color, extent

» Positioning

- Relative to landmarks reported by “dlib”
- If “dlib” fails, then relative to detected eyes from good FRVT FR algorithms

1
ORIGINAL
IMAGE



2
WIDE,
HIGH
COVERAGE



3
WIDE,
MEDIUM
COVERAGE



4
WIDE,
LOW
COVERAGE



5
ROUND,
HIGH
COVERAGE



FRVT Leaderboard (all without masks)

<https://pages.nist.gov/frvt/html/frvt11.html>

Developer	VISA Photos FNMR @ FMR ≤ 0.000001	MUGSHOT Photos FNMR @ FMR ≤ 0.00001	MUGSHOT Photos FNMR @ FMR ≤ 0.00001 DT>=12 YRS	VISABORDER Photos FNMR@FMR ≤0.000001	BORDER Photos FNMR @ FMR = 0.000001	WILD Photos FNMR@ FMR ≤ 0.00001	CHILD EXP Photos FNMR@ FMR ≤ 0.01
visionlabs-008	0.0036 ⁽⁴⁾	0.0031 ⁽⁶⁾	0.0040 ⁽⁴⁾	0.0045 ⁽³⁾	0.0079 ⁽¹⁾	0.0308 ⁽¹⁰⁾	-
ntechlab-008	0.0061 ⁽¹⁰⁾	0.0056 ⁽¹⁷⁾	0.0108 ⁽²⁹⁾	0.0042 ⁽¹⁾	0.0080 ⁽²⁾	0.0312 ⁽²⁰⁾	-
deepglint-002	0.0027 ⁽²⁾	0.0032 ⁽⁷⁾	0.0033 ⁽²⁾	0.0043 ⁽²⁾	0.0084 ⁽³⁾	0.0301 ⁽¹⁾	0.3422 ⁽⁶⁾
dahua-004	0.0058 ⁽⁹⁾	0.0036 ⁽⁸⁾	0.0048 ⁽⁸⁾	0.0051 ⁽⁵⁾	0.0086 ⁽⁴⁾	0.0304 ⁽⁵⁾	-
vocord-008	0.0038 ⁽⁵⁾	0.0042 ⁽¹¹⁾	0.0055 ⁽¹²⁾	0.0045 ⁽⁴⁾	0.0086 ⁽⁵⁾	0.0310 ⁽¹⁴⁾	-
cuhkee-001	0.0045 ⁽⁶⁾	0.0031 ⁽⁵⁾	0.0046 ⁽⁷⁾	0.0051 ⁽⁷⁾	0.0095 ⁽⁶⁾	0.1524 ⁽¹⁰²⁾	-
sensetime-003	0.0027 ⁽³⁾	0.0027 ⁽¹⁾	0.0027 ⁽¹⁾	0.0051 ⁽⁶⁾	0.0100 ⁽⁷⁾	0.0355 ⁽⁴⁵⁾	0.3683 ⁽⁷⁾
alleyes-000	0.0090 ⁽²¹⁾	0.0055 ⁽¹⁵⁾	0.0087 ⁽²¹⁾	0.0068 ⁽¹⁰⁾	0.0105 ⁽⁸⁾	0.0306 ⁽⁸⁾	-
tech5-004	0.0234 ⁽⁷²⁾	0.0086 ⁽⁴²⁾	0.0162 ⁽⁵³⁾	0.0065 ⁽⁹⁾	0.0112 ⁽⁹⁾	0.0311 ⁽¹⁷⁾	-
yitu-003	0.0026 ⁽¹⁾	0.0066 ⁽²⁵⁾	0.0085 ⁽¹⁷⁾	0.0064 ⁽⁸⁾	0.0114 ⁽¹⁰⁾	0.0360 ⁽⁴⁹⁾	-

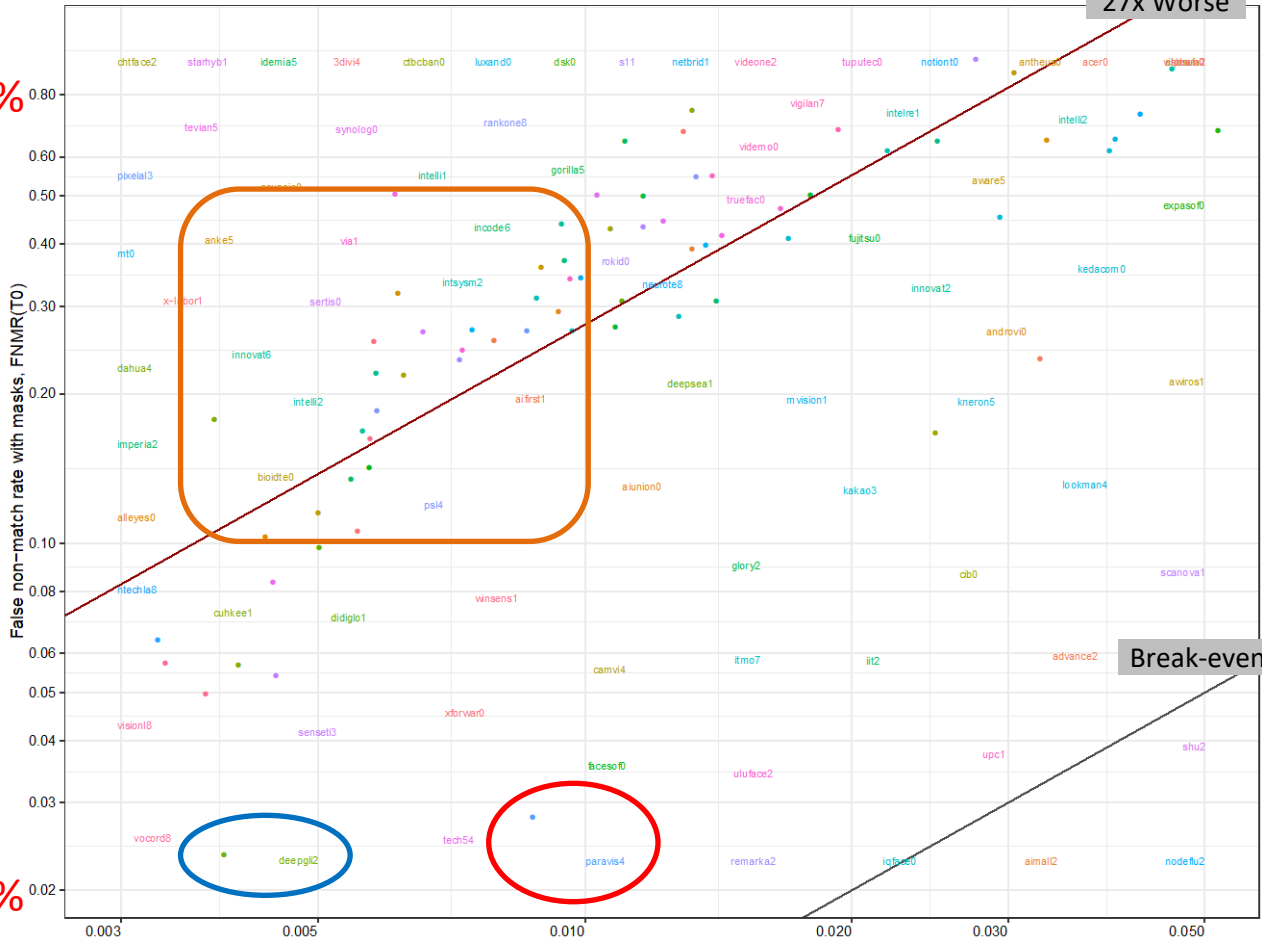
Accuracy with and without masks

False Negative Rate With Masks

Impact of medium wide lightblue masks
The lower line is $y = x$; the upper line is $y = 27.6x$

80%

27x Worse



0.3%

False non-match rate without masks, FNMR(T0) with FMR(T0) = 0.00010

False Negative Rate Without Masks

5%

But... further challenges



Some algorithms may be usable

Most pre-pandemic algorithms do not tolerate masks

Algorithm	VISABORDER Photos FNMR @ FMR ≤ 0.00001 (NOT MASKED)	VISABORDER Photos FNMR@FMR ≤ 0.00001 (MASKED PROBE) lightblue, wide, medium coverage
deepglint-002	0.0039 ⁽⁹⁾	0.0237 ⁽¹⁾
paravision-004	0.0088 ⁽⁴⁸⁾	0.0281 ⁽²⁾
visionlabs-009	0.0028 ⁽¹⁾	0.0355 ⁽³⁾
iqface-002	0.0086 ⁽⁴⁶⁾	0.0445 ⁽⁴⁾
pensees-001	0.0106 ⁽⁶⁰⁾	0.0461 ⁽⁵⁾
vocord-008	0.0038 ⁽⁷⁾	0.0500 ⁽⁶⁾
idemia-006	0.0048 ⁽¹⁷⁾	0.0539 ⁽⁷⁾
rankone-008	0.0134 ⁽⁵²⁾	0.5470 ⁽⁵⁸⁾
videmo-000	0.0140 ⁽⁵⁴⁾	0.5509 ⁽⁵⁹⁾
scanovate-001	0.2403 ⁽⁸⁰⁾	0.5973 ⁽⁶⁰⁾
intelresearch-001	0.0220 ⁽⁶¹⁾	0.6184 ⁽⁶¹⁾
kedacom-000	0.0391 ⁽⁷¹⁾	0.6188 ⁽⁶²⁾
innovativetechnologyltd-002	0.0251 ⁽⁶⁴⁾	0.6454 ⁽⁶³⁾
idemia-005	0.0111 ⁽⁴⁴⁾	0.6469 ⁽⁶⁴⁾

Failure to verify rate rises from 0.4% to 2.4%

Failure to verify rate rises from 1% to 65%



Demographic Effects

FR accuracy varies by population

Landscape

- Race? Sex? Age? What else?
- Algorithms, cameras?
- 1:1 vs. 1:N
- False positives? Or Negatives?

NIST tests and results

- Criminal investigation
- Clustering media

Scope of NIST demographics work

» Algorithms

- 187 algorithms, 99 developers
- Mostly commercial, some universities
- Prototypes from R&D labs

» Modes

- One to one verification (DHS, DoS)
- One to many identification (mugshots)

» Metrics:

- False positives
- False negatives
- Failure to enroll

» Relevance to applications

» 18.3 million cooperative photos of 8.5 million people

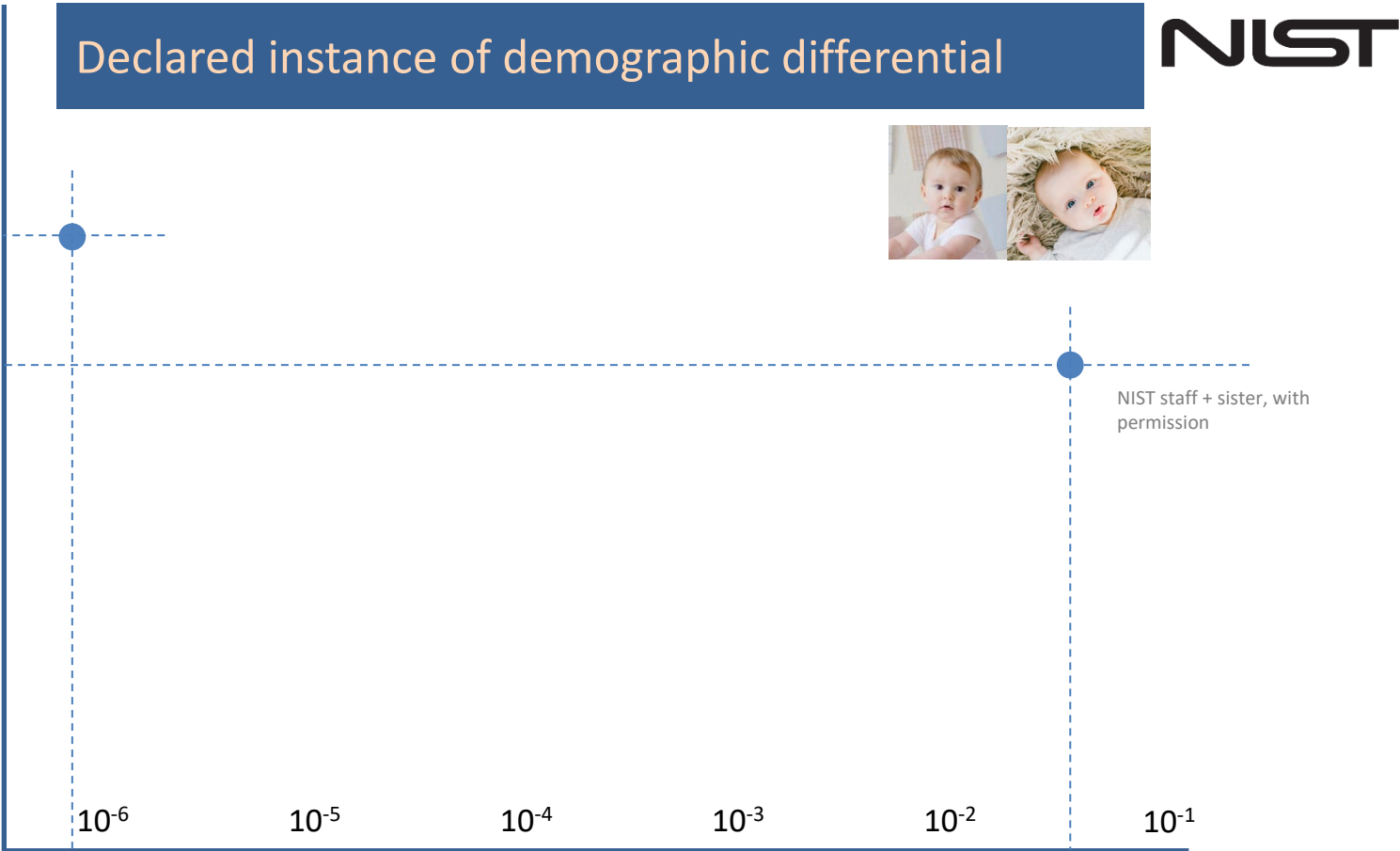
- **DHS/CIS Application Photos**
 - High quality
 - Race: 24 countries, 7 regions
 - Sex: M, F only
 - Age groups: [12-20], [20-35], [35-50], [50-65], [65-99].
- **DHS/CBP Entry Photos**
 - Mediocre quality
 - Compare with CIS photos
- **DOS Visa photos**
 - Age
- **FBI mugshots**
 - Sex: M, F, only
 - Age groups: Adults above or below 45.
 - Race: Asian, Black, White, Native American

Declared instance of demographic differential

FNMR
False non-match rate
Proportion of genuine comparisons producing score below threshold, T .
See ISO/IEC 19795-1



Log-scale is typical to show small numbers.



NIST staff + sister, with permission

Log-scale is often required because low FMR values are operationally relevant.

FMR False match rate
Proportion of impostor comparisons searches yielding any candidates at or above threshold, T .⁴⁶

Declared instance of demographic differential

FNMR
False non-match rate

Proportion of genuine comparisons producing score below threshold, T .

See ISO/IEC 19795-1

Log-scale is typical to show small numbers.

10^{-6}

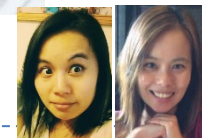
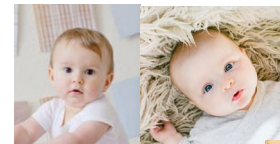
10^{-5}

10^{-4}

10^{-3}

10^{-2}

10^{-1}



NIST staff + sister, with permission

Apple Face. ID claims FMR $\sim 1:1\,000\,000$
<https://support.apple.com/en-us/HT208108>

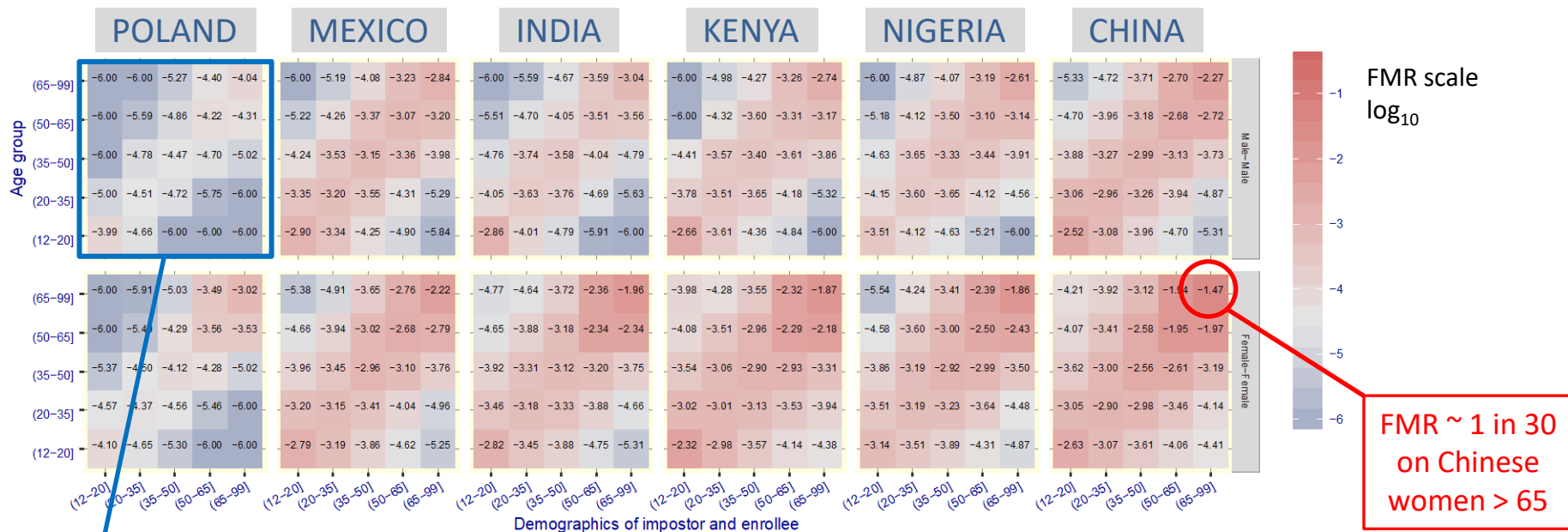
“The statistical probability is different for twins and siblings that look like you and among children under the age of 13, because their distinct facial features may not have fully developed.”

Log-scale is often required because low FMR values are operationally relevant.

FMR False match rate

Proportion of impostor comparisons searches yielding any candidates at or above threshold, T .⁴⁷

Cross-age false match rates in six countries, male x male, and female x female



FMR on white males are below 1 in 30 000

Algorithm Imperial-002 with T = 1.381120
 Nominal FMR = 0.00003
 Dataset = Frontal cf. passports

Source: NIST IR 8280, 2019-12

Thinking through consequences: Three applications

1. Dispensing drugs

- » Non-repudiation
- » 1:1
- » Volume: 100s per day
- » Transactions are almost always mated
 - Prob(Impostor) is LOW
- » False negative → Inconvenience
- » False positive → Prescription drug fraud

- » Who is harmed by demographic differential in FP?
 - Some pharmacists

2. Boarding a plane

- » Facilitation of recording immigration exit vs. Access Control
- » 1:N
- » Volume: 100s per flight
- » Transactions are almost always mated
 - Prob (Impostor) is LOW
- » False negative → Paper boarding with airline staff
- » False positive → Stowaway
 - but manifest exists, and legitimate customer may board also so “low” consequences
- » Who is harmed by demographic differential in FP?
 - Airline.

3. Watchlist

- » Soccer stadium. Counter-terrorism. Compulsive gamblers
- » 1:N
- » Volume: 10s of thousands per day
- » Transactions are almost always non-mated
 - Prob (Genuine) is LOW
- » False negative → Undetected “bad guy”
- » False positive → Incorrect enforcement action ... civil liberties

- » Who is harmed by demographic differentials in FP?
 - Bystanders

- » Leading contemporary algorithms
 - Are very accurate
 - Increasingly tolerate poor image quality
 - Generally distribute errors inequitably across demographics

- » Leading contemporary algorithms
 - Are very accurate
 - Increasingly tolerate poor image quality
 - Generally distribute errors inequitably across demographics

- » False positive differentials much larger than false negative differentials
 - More false positives in Asian and African faces
 - More false positives in women
 - More false positives in the old and very young

- » Leading contemporary algorithms
 - Are very accurate
 - Increasingly tolerate poor image quality
 - Generally distribute errors inequitably across demographics

- » False positive differentials much larger than false negative differentials
 - More false positives in Asian and African faces
 - More false positives in women
 - More false positives in the old and very young

- » One-to-many algorithms don't necessarily behave like one-to-one
 - Some one-to-many effect a stabilization of the impostor distribution

- » Leading contemporary algorithms
 - Are very accurate
 - Increasingly tolerate poor image quality
 - Generally distribute errors inequitably across demographics

- » False positive differentials much larger than false negative differentials
 - More false positives in Asian and African faces
 - More false positives in women
 - More false positives in the old and very young

- » One-to-many algorithms don't necessarily behave like one-to-one
 - Some one-to-many effect a stabilization of the impostor distribution

- » Algorithm matters
 - Accuracy
 - Demographic sensitivity
 - Know-your-algorithm
 - Traceability to (NIST) tests is not easy

- » Leading contemporary algorithms
 - Are very accurate
 - Increasingly tolerate poor image quality
 - Generally distribute errors inequitably across demographics

- » False positive differentials much larger than false negative differentials
 - More false positives in Asian and African faces
 - More false positives in women
 - More false positives in the old and very young

- » One-to-many algorithms don't necessarily behave like one-to-one
 - Some one-to-many effect a stabilization of the impostor distribution

- » Algorithm matters
 - Accuracy
 - Demographic sensitivity
 - Know-your-algorithm
 - Traceability to (NIST) tests is not easy

- » Application matters
 - Error impact can be grave or inconsequential.

- » Leading contemporary algorithms
 - Are very accurate
 - Increasingly tolerate poor image quality
 - Generally distribute errors inequitably across demographics
- » False positive differentials much larger than false negative differentials
 - More false positives in Asian and African faces
 - More false positives in women
 - More false positives in the old and very young
- » One-to-many algorithms don't necessarily behave like one-to-one
 - Some one-to-many effect a stabilization of the impostor distribution
- » Algorithm matters
 - Accuracy
 - Demographic sensitivity
 - Know-your-algorithm
 - Traceability to (NIST) tests is not easy

- » Application matters
 - Error impact can be grave or inconsequential.
- » Incomplete reporting in the press and academia
 - Confusion of face “analysis” with “recognition”
 - Don't identify which component is at fault
 - Missing reports on false positives
 - Differentiate false positives from false negatives

Twins: The Forgotten Demographic



Source: Twins Day Ohio collected by Notre Dame

Same person or not?



	Identical	Fraternal
How	Monozygotic	Dizygotic
Proportion of individuals that are a twin	0.9%	3.1%
Same-sex	100%	50% in theory 58% actually
TR gain since 1980	x1.5 since 1980	x1.9 since
Demographics	~ constant with age, geography	varies with mothers age, order, geography

Twins, triplets ... constituted 140,000 out of 4M births in 2015
https://www.cdc.gov/nchs/data/nvsr/nvsr66/nvsr66_01.pdf

Scenario: Identical Twins

Probe is an identical twin



Gallery Size: 1.6 million

Algorithm	Rank of sibling	Score	FPIR
Microsoft	1	0.78	0.0007
NEC	1	0.77	0.0010
Idemia	1	3066	0.0007

Almost all algorithms give high scores

Candidate List



Rank 1

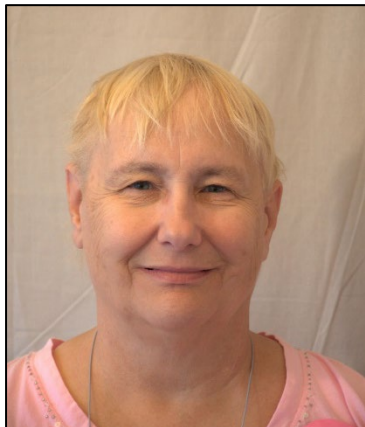


...

...

Scenario: Fraternal Twins

Probe is a fraternal twin



Gallery Size: 1.6 million

Algorithm	Rank of sibling	Score	FPIR
Microsoft	1	0.18	0.878
NEC	1	0.64	0.986
Idemia	11	670	0.909

Candidate List



Rank 1
(NEC/Microsoft)



...



Rank 11
(Idemia)

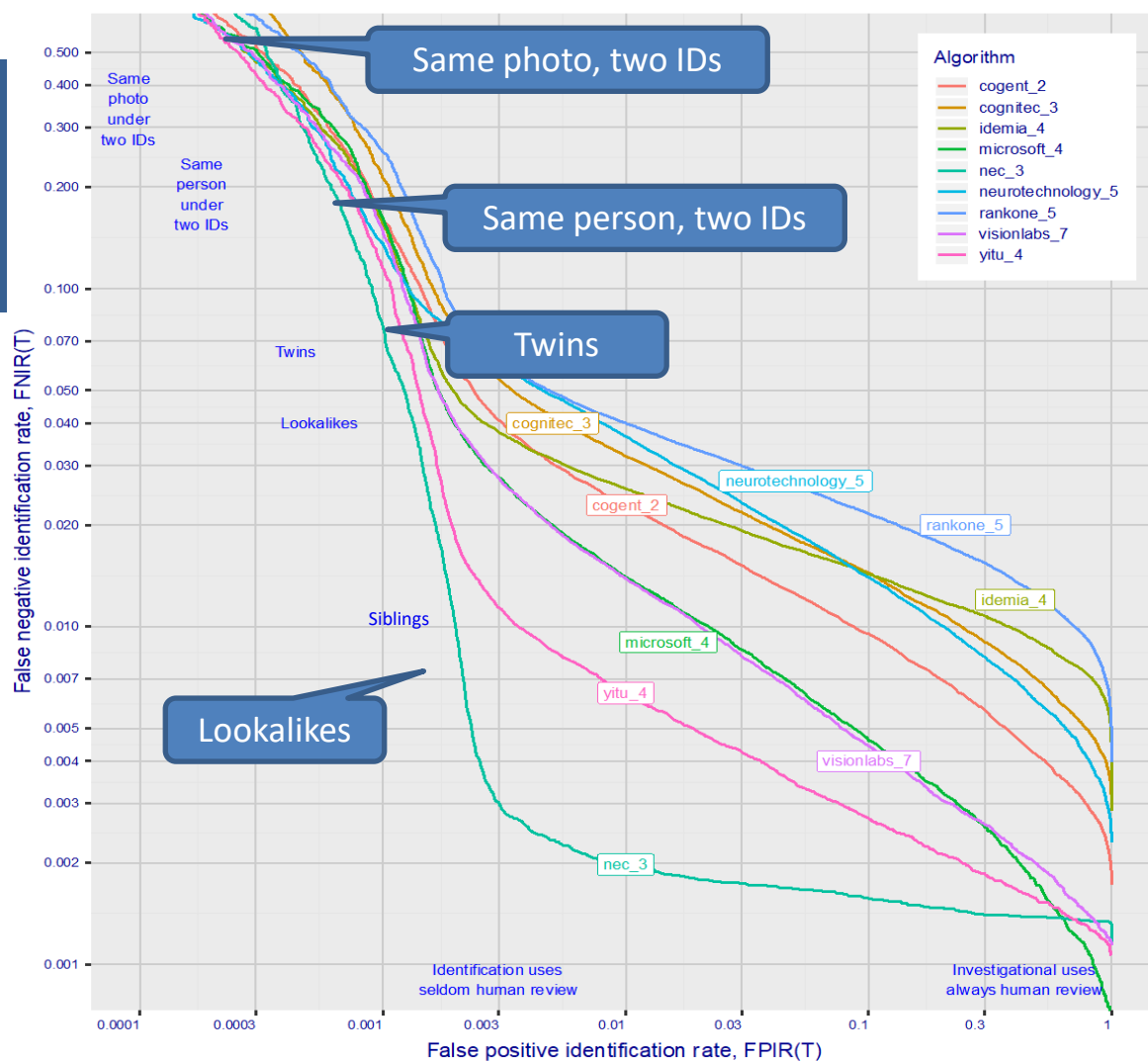


...

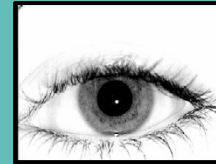
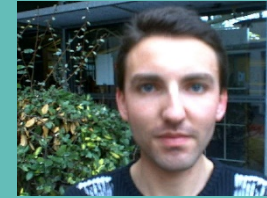
Face Recognition at National Scale

In a “closed” population (town, country):

- Low false positive rates cannot be achieved due to familial relationships
- Not expected with 10 fingerprints, and iris recognition



Why Face? Versus Fingerprint, Iris.



Source: <http://biometrics.it-sudparis.eu/english/index.php?menu=datasample>

Modality selection

Modality	Image appearance standards	Availability (Ease of capture)	Permanence (ageing)	Uniqueness	Demographic problems	Twins	Retained reference images	Social acceptance
Face	Yes, compliance is difficult and not necessary	Fast Non-contact Socially accepted	Lower Low in children	Lower	Strong false positives in twins, families, same ethnicities, same sex, age	FMR → 1 identical twins FMR high also in fraternal	Social media, gov databases, (passport, drivers license)	Highest: Global ICAO passport
Finger Contact	Yes	Single fastest Four fast Ten slow (for gov use)	High Possibility of environmental damage	High Very high 10 fingers	No More false negatives in the elderly, very young, depends on sensor	FMR → 0	Legacy gov databases	Lower: Local cultural
Finger Contact-less	No: Interoperability problems with contact	Fast Four fingers for physical access control	High	High	No	FMR → 0	Yes, but only contact fingerprints	Higher: For PACS
Iris	Partial Guidance yes	Slower, optical tradeoffs. Capture both simultaneously	High, possibility of disease	High Very high two irides	No False negatives in elderly	FMR → 0	Few	Lower



- Nuanced discussion around many of these entries
- There are applications where property is not relevant

ONGOING BENCHMARKS

- 1. FRVT 1:1**
Core Biometric Operation
- 2. FRVT 1:N**
Search Performance
- 3. FRVT Morph**
Morphed Photo Detection
- 4. FRVT Quality**
Automated Quality Assessment

CURRENT PRODUCTS

Part 1: Performance of 1:1 Verification Algorithms	Part 2: Performance of 1:N Identification Algorithms	Part 3: Demographic Effects in Face Recognition	Part 4: Performance of Morph Detection Algorithms	Part 5: Performance of Image Quality Assessment Algorithms	Part 6: Performance of Face Recognition with Face Masks	Part 7: Performance of Face Recognition on Twins
<p>Last: 2020-08-25 Next: 2020-07</p>	<p>Last: 2020-03-27 Next: 2020-08</p>	<p>Last: 2019-12-19 Next: 2020-09</p>	<p>Last: 2020-07-24 Next: 2020-09</p>	<p>Last: 2020-07-27 Next: 2020-09 est.</p>	<p>Last: 2020-07-27 Next: 2020-08 est.</p>	<p>Last: Next: TBD</p>

Impeding accuracy

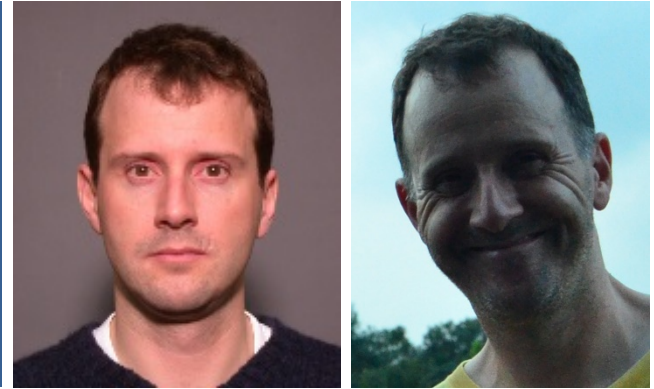
- » Ageing
- » Twins
- » Demographics differentials
 - False positives WORSE THAN false negatives
- » Poor quality images
 - Pose
 - Illumination
 - Resolution
 - Occlusion (face masks)
 - Cropping
 - Distortion
- » Lack of capture standards

Impeding security

- » Morph attack detection
- » Presentation attack detection
- » Tampering
- » Fakes

ISO/IEC 24358

FASTER, BETTER, FACE-AWARE CAPTURE (QUALITY MATTERS!)



Images from presenter

Problems:

- a) Non-frontal faces
- b) No-faces, multiple-faces
- c) Over-, under-exposure
- d) Human review errors
- e) Morphing
- f) Inadequate presentation attack detection

Potential Solutions:

- a) Face pose detector
- b) Face detectors
- c) 12 bits or closed-loop control
- d) Higher resolution, better compression, 3D
- e) Crypto for tamper-proofing

NIST IFPC Conference: October 27-29.



IFPC 2020 - Tuesday Oct 27		IFPC 2020 - Wednesday Oct 28		IFPC 2020 - Thursday Oct 29	
	07:20 Welcome		07:00 Welcome		07:00 Welcome
11	07:30 Arun Vemury , DHS Science + Technology Directorate (US): <i>Welcome + DHS context</i>	21	07:10 Lars Ericson , IARPA (US): <i>Overview of the IARPA efforts on face recognition</i>	31	07:10 Rebecca Heyer , DSTG (AU): <i>Face recognition in Australia</i>
12	07:40 Istvan Szilard Racz , EU-LISA: <i>European Entry-Exit System</i>	22	07:40 Stergios Papadakis , Johns Hopkins Applied Physics Lab (US): <i>Results from the Odin program on presentation attack detection</i>	32	07:40 Martins Bruveris , Onfido (UK): <i>Reducing geographic performance differentials for face recognition</i>
13	08:10 Anna Stratmann , BSI (DE): <i>Biometric processes of the Entry Exit System</i>	23	08:10 Marta Gomez-Barrero , Hochschule Ansbach (DE): <i>Presentation attack detection and unknown attacks</i>	33	08:10 Mosalam Ebrahimi , Trueface AI (US): <i>A bias mitigation strategy: overcoming the problem of overly confident false matches</i>
14	08:40 Patrick Grother , NIST (US): <i>Measurement of face recognition performance for Entry-Exit</i>	24	08:40 Christian Rathgeb , Hochschule Darmstadt (DE): <i>Impact of facial beautification on face recognition: From plastic surgery to makeup presentation attacks</i>	34	08:40 Jacqueline Cavazos , UT Dallas (US): <i>Accuracy comparison across face recognition algorithms: Where are we on measuring race bias?</i>
	09:10 <i>Break 15 mins</i>		09:10 <i>Break 15 mins</i>		09:10 <i>Break 15 mins</i>
15	09:25 Arun Ross , Michigan State University (US): <i>Look-alike disambiguation in face recognition</i>	25	09:25 Stéphane Gentic , Idemia (FR): <i>Synthetic faces: Are they new identities and can they be used in evaluation?</i>	35	09:25 John Howard & Yevgeniy Sirotnin , SAIC (US): <i>Revisiting the Fitzpatrick Scale and Face Photo-based Estimates of Skin Phenotypes</i>
16	09:55 P. Jonathon Phillips , NIST (US): <i>Item response theory for designing calibrated face ability tests</i>	26	09:55 Mei Ngan , NIST (US): <i>Face morphing - threats, technology, what's next</i>	36	09:55 Michael Thieme , Novetta (US): <i>AI performance assessment standardization in SC 42 – implications for biometrics</i>
17	10:25 Laura Rabbitt & Yevgeniy Sirotnin , SAIC (US): <i>Human-Algorithm Teaming in Face Recognition</i>	27	10:25 Christoph Busch , NTNU/Hochschule Darmstadt (NO/DE): <i>Face morphing attack detection in the iMARS project</i>	37	10:25 Johanna Morley , Metropolitan Police (UK): <i>Testing of demographic effects in an operational live facial recognition from video system</i>
18	10:55 Carina A. Hahn , NIST (US): <i>The effectiveness of fusion in face recognition</i>	28	10:55 Kiran Raja , NTNU/MOBAl (NO): <i>Morphing Attack Detection - obstacles for research to deployment</i>	38	10:55 Brendan Klare , Rank One Computing (US): <i>Efficiency considerations for face recognition algorithms</i>
19	11:25 Amy N. Yates , NIST (US): <i>Perceptual face abilities of face examiners for varying tasks</i>	29	11:25 Chen Liu, Zander Blasingame , Clarkson U., David Doermann , U. at Buffalo, Jeremy Dawson , West Virginia U. (US): <i>Center for Identification Technology Research (CITeR) Morph Attack Detection and Mitigation Projects</i>	39	11:25 Bhargav Avasarala , Paravision (US): <i>Challenges and considerations for masked face recognition</i>
1a	11:55 John Howard & Yevgeniy Sirotnin , SAIC (US): <i>Quantifying Race and Gender Effects in Face versus Iris Algorithms</i>	2a	11:55 Pawel Drozdowski Hochschule Darmstadt (DE): <i>Workload reduction in FR identification with morphing</i>	3a	11:55 Tony Mansfield , NPL (UK): <i>The new ISO/IEC 19795-1 biometric performance testing and reporting standard</i>
1b	12:25 Patrick Grother , NIST (US): <i>Now under development: ISO/IEC 29794-5 face image quality standard ISO/IEC 24358 face-aware capture specifications</i>	2b	12:25 Mei Ngan , NIST (US): <i>Evaluation of face recognition accuracy for subjects potentially wearing face masks</i>	3b	12:25
	12:55 <i>Close</i>		12:55 <i>Close</i>		12:55 <i>Close</i>

THANKS

PATRICK.GROTHER@NIST.GOV

FRVT@NIST.GOV

Facial Recognition Performance and Its Measurement



PRESENTED BY:

Patrick Grother

National Institute of Standards and Technology

MODERATED BY:

Stephen Redifer

2020-09-24



HDIAC

Homeland Defense & Security
Information Analysis Center

DISTRIBUTION STATEMENT A. Approved for public release; distribution is unlimited.

HDIAC is sponsored by the Defense Technical Information Center (DTIC). Any opinions, findings and conclusions or recommendations expressed in this material are those of the author(s) and do not necessarily reflect the views of the Defense Technical Information Center.

info@hdiac.org
<https://www.hdiac.org>