# NIST Interagency Report 8429
## Summarizing Demographic Differentials

Patrick Grother
NIST

IFPC
November 17, 2022

**Why?**
Demographics *do* have an effect

**What?**
Quantifying the problem

**Role?**
As a target for designers to optimize

**Validate how?**
Exercise them (in FRVT…)

**Standardize where?**
ISO/IEC 19795-10 Demographics

**For use in (for example)**
ISO/IEC 9868 EU Regulation

# Quoting Georgetown's Report: "The Perpetual Line-up"

**[Bias Exists]**
- The most prominent study [Klare et al.] found that several leading algorithms performed worse on African Americans, women, and young adults than on Caucasians, men, and older people, respectively.[216]

**[Consequence]**
- If the suspect is African American rather than Caucasian, the system is more likely to erroneously fail to identify the right person, potentially causing innocent people to be bumped up the list—and possibly even investigated

**[Awareness]**
- "Q: Is the Booking Photo Comparison System biased against minorities[?]"
- "A: No… it does not see race, sex, orientation or age. The software is matching distance and patterns only, not skin color, age or sex of an individual."- Frequently Asked Questions, Seattle Police Department

**[No Bias Tests]**
- There is no independent testing regime for racially biased error rates … two major face recognition companies admitted that they did not run these tests

**[Priors]**
- Racial bias intrinsic to an algorithm maybe compounded by outside factors. African Americans are disproportionately likely to come into contact with—and be arrested by—law enforcement.[218]

Clare Garvie, Alvaro M. Bedoya, Jonathan Frankle
*The Perpetual Line-up Unregulated Police Face Recognition In America*
Georgetown Law Center on Privacy and Technology
October 18, 2016  https://www.perpetuallineup.org/

2

# NIST Interagency Report 8280 - Dec 2019

- Distinguish between False Negatives and False Positives
- Distinguish between 1:1 and 1:N
- Consequences of differentials are application dependent
- Effects are algorithm dependent $\longrightarrow$ know your algorithm $\longrightarrow$ know your system
- Mitigation guidance

---

- ΔFNMR small cooperative images
- ΔFMR massive even in cooperative images
- Higher FNMR and FMR in women
- Higher FMR in East Asia and Africa
- Some Chinese algorithm give higher FMR in Europe
- Some 1:N algorithms effect low  ΔFPIR

# Effect of age:

# False Negative Identification Rates aka "Miss Rates"

Algorithm: Canon-001 (2021-10-27)
Images: Airport immigration photos

False Negative Identification "Miss" Rates. N = 1.6 million
The threshold set to limit false positive outcomes to 1 in
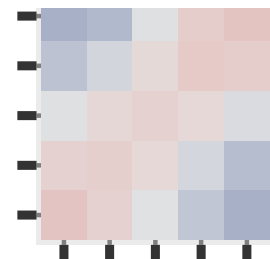1000 searches (FPIR = 0.001) for men age 30-45.

| AGE | Female | Male |
|---|---|---|
| (60:99] | 0.026 | 0.013 |
| (45:60] | 0.022 | 0.009 |
| (30:45] | 0.025 | 0.010 |
| (21:30] | 0.035 | 0.015 |
| (18:21] | 0.060 | 0.046 |
| (15:18] | 0.088 | 0.128 |
| (12:15] | 0.115 | 0.226 |
|  | 10 | 10 |

# 1:N False Positive Rates by Sex and Age

Algorithm: canon_001, Dataset: Border – Border, N = 1600000
Threshold: 1.442880 for FPIR(T, 30–45, Male) = 0.001
Text encodes FPIR, Color encodes log(FPIR)

−4  −3  −2  −1
FMR 1 in 10000   FMR 1 in 100

**Age group of person in non−mate probe**

| Female | Male |
|--------|------|
| (60:99]  0.0124 | 0.0006 |
| (45:60]  0.0073 | 0.0004 |
| (30:45]  0.0088 | 0.0010  FPIR = 1 IN 1000 |
| (21:30]  0.0173 | 0.0026 |
| (18:21]  0.0248 | 0.0047 |
| (15:18]  0.0279  FPIR = 1 IN 35 | 0.0069 |
| (12:15]  0.0305 | 0.0107 |

**Sex of person in non−mate probe**

5

# Cross age false match rates



Higher FMR
- Young
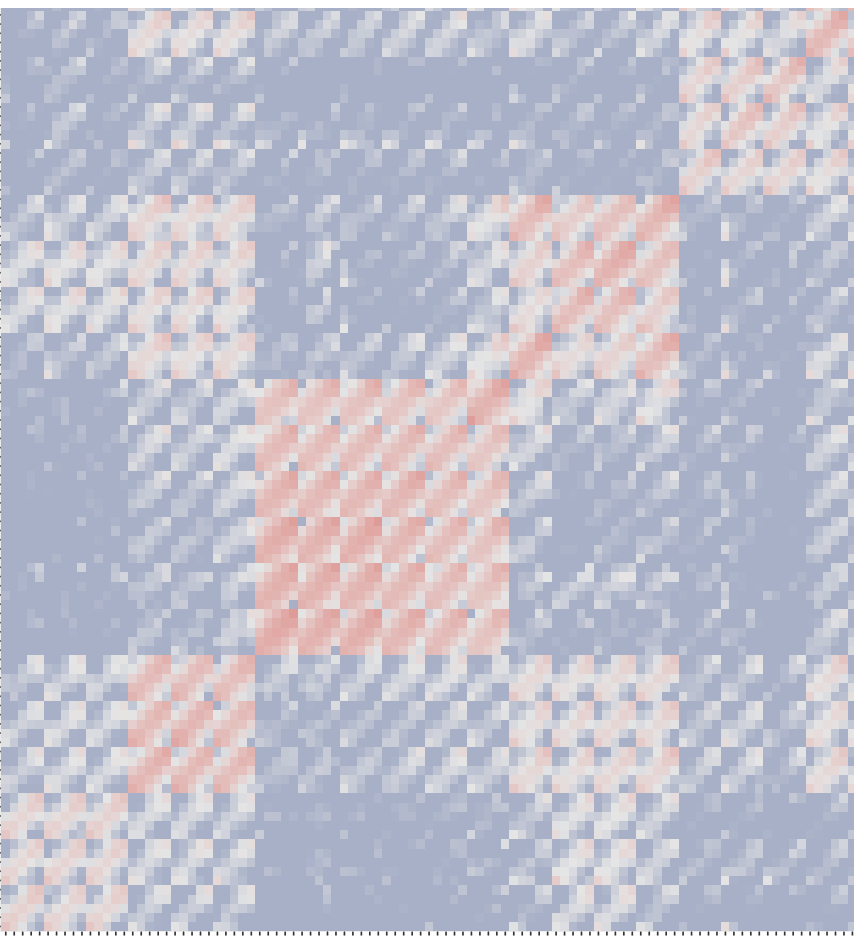- Old

Lower FMR
- Cross-age

Algorithm: dahua_003  Threshold: 6430.000000  Dataset: Application
Nominal FMR: 0.000030  log10 FMR

6

# Cross country-of-birth and age false match rates

Adding women born in:
- Russia
- Ukraine

# Cross country-of-birth and age false match rates

Adding people born in:

- El Salvador
- Mexico
- Nicaragua



Algorithm: dahua_003  Threshold: 6430.000000  Dataset: Application
Nominal FMR: 0.000030  log10 FMR

# Cross country-of-birth and age false match rates

Adding women born in:
- Nigeria
- Liberia
- Ghana

Algorithm: dahua_003  Threshold: 6430.000000  Dataset: Application
Nominal FMR: 0.000030  log10 FMR

# Cross country-of-birth and age false match rates

Adding women in
- 22 countries
- 7 regions
  - E. Europe
  - C. America
  - W. Africa
  - Caribbean
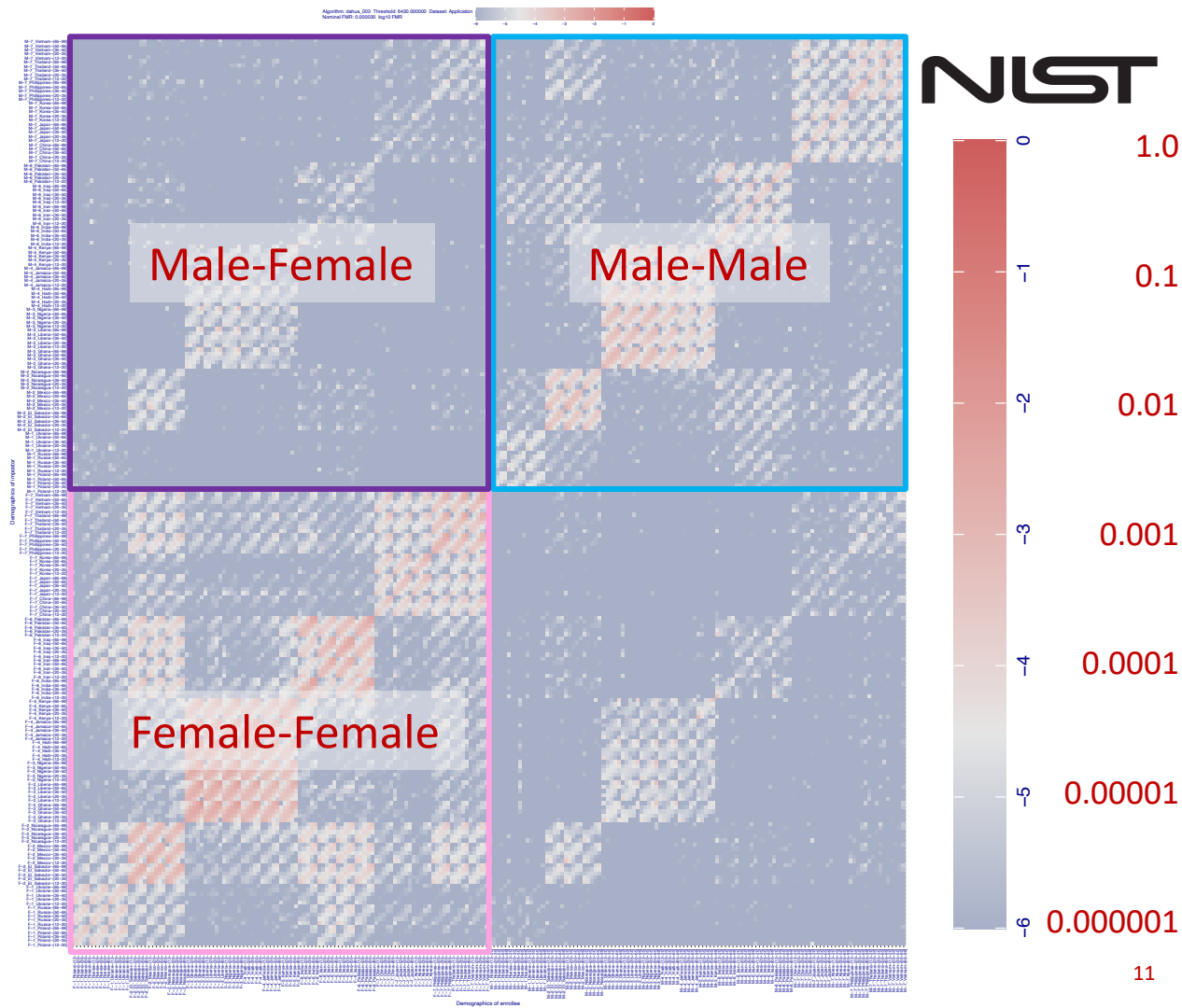  - E. Africa
  - S. Asia
  - E. Asia



Algorithm: dahua_003  Threshold: 6430.000000  Dataset: Application
Nominal FMR: 0.000030  log10 FMR

# Cross country-of-birth, age, and sex false match rates

Adding:
- Males

# How false positives affect 1:N applications

Gallery composition:

1. Six demographic groups
2. Equally balanced – all have same number of people

Gallery composition:

1. Six demographic groups
2. Now imbalanced - the usual case
3. Number of people in group i is $n_i$

4. Gallery size is N = $\sum n_i$

NIST



Brondby fans scuffle with police during a match between the Copenhagen and Brondby soccer teams at Copenhagen's Telia Parken stadium in 2017.

*Lars Ronbog/FrontzoneSport via Getty Images*

Num. enrolled ~ 50
Num. searches ~ 21000

Once the men's chant is over, the group moves toward the stadium's entrance, where the men — along with 21,000 other fans — are asked to remove masks, hats and glasses so a computer can scan their faces. The scans will be compared against a list of roughly 50 banned troublemakers and will be used to determine whether the spectators will be allowed in.

No one is stopped on this day. But since the system's launch in July, it has caught four people on the blacklist, who were then turned over to police.

https://www.npr.org/2019/10/21/770280447/a-soccer-team-in-denmark-is-using-facial-recognition-to-stop-unruly-fans?sc=tw&t=1572190088133

Probe composition over some time period

1. Say a total of 21000 searches
2. Almost all non-mates
3. Again imbalanced
4. Number of people in group i is $p_i$

5. Potentially gallery and probe composition differ

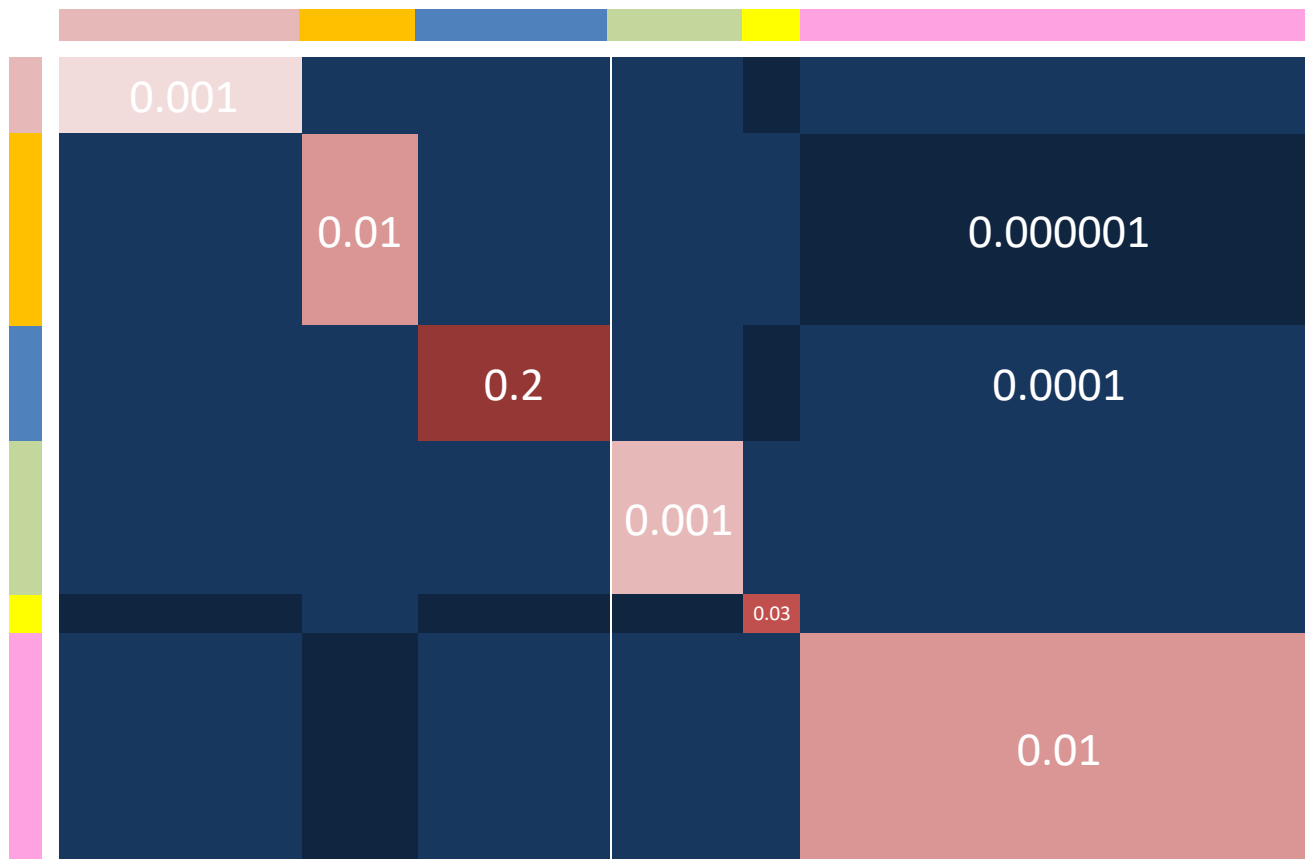# How false positives affect 1:N applications

Number of expected false positives

- NFP   = N  FMR(T)  P

where
- N       =             gallery size
- P       =             number of non-mated searches
- FMR   =             monolithic 1:1 comparison false match rate
-                       assuming FMR doesn't depend on demographics

# How false positives affect 1:N applications



$$\text{NFP} = \sum \sum p_j FMR_{ji} n_i$$
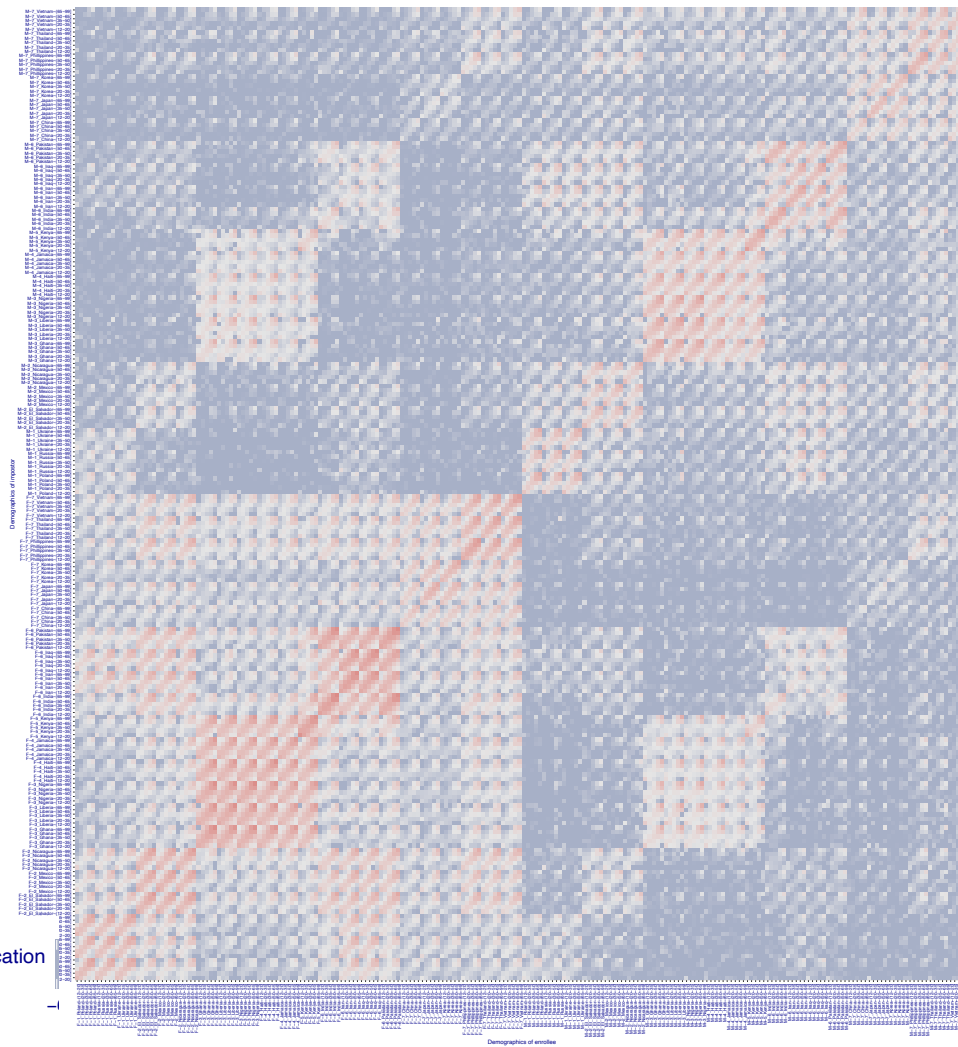
In this toy case:

NFP is dominated by the high FMR group which has 20x higher FMR than any other.
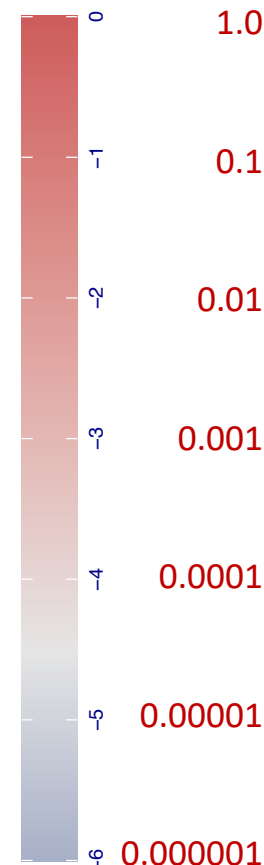
See NISTIR 8429 Annex B

# Cross country-of-birth, age, and sex false match rates

Adding:
- Males

# Prior publication demographic consequences of FMR differentials on one-to-many search

**Quantifying the Extent to Which Race and Gender Features Determine Identity in Commercial Face Recognition Algorithms**

John J. Howard, Yevgeniy B. Sirotin, Jerry L. Tipton, A. Vemury

Published 15 October 2020

Also see the older literature on (binomial) models of 1:N accuracy with heterogeneous error rates.

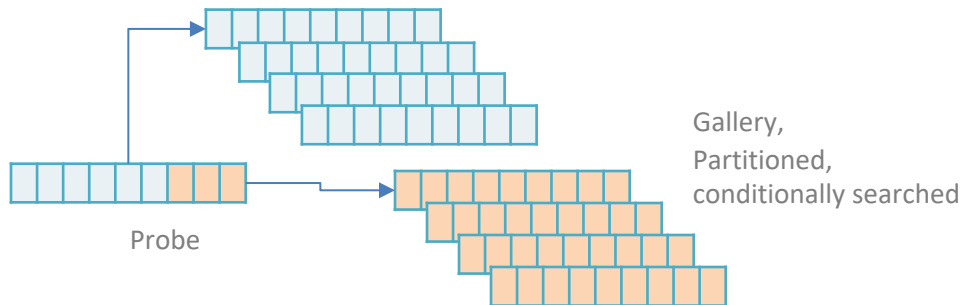https://mdtf.org/publications/TPS-Features.pdf

Caveat esp. for reviewers coming back to this deck later

- Many developers implement 1:N search as N 1:1 comparisons
- **But some do not**
    - **The enrollment database is not just N separate templates**
    - **It could be a tree, or a dictionary, or some exotic data structure**
- Some developers field both types of algorithms

- This has beneficial consequences for:
    - False positive rates
    - How false positive rates grow when N goes up
    - Demographic dependencies
    - Speed
- This has complexity
    - Deleting somebody from a database may not be a simple operation

Gallery, Partitioned, conditionally searched

Probe



Yu A. Malkov, D. A. Yashunin, **Efficient and Robust Approximate Nearest Neighbor Search Using Hierarchical Navigable Small World Graphs.** IEEE Transactions on Pattern Analysis and Machine Intelligence, Volume 42 No. 4 April 2020 pp. 824–836 https://doi.org/10.1109/TPAMI.2018.288947
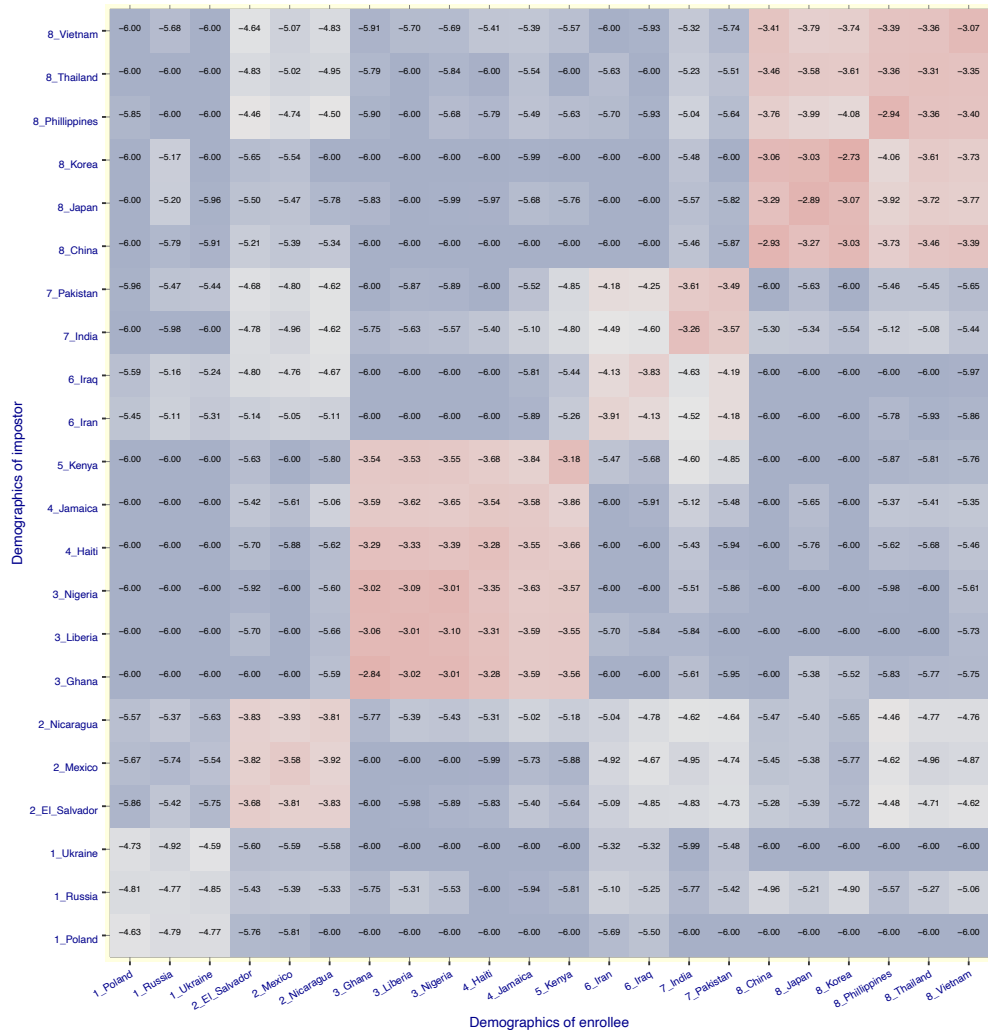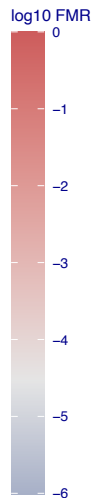
Fig. 1. Illustration of the Hierarchical NSW idea. The search starts from an element from the top layer (shown red). Red arrows show direction of the greedy algorithm from the entry point to the query (shown green).

3

# Cross country FMR: (India)



1. FMR EU ~ 1:33000
2. FMR Nigeria 1:1000
3. FMR Korea 1:500
4. Relevance to 1:N

Magnitude matters: Age x Age for six countries

# Candidate Measures

**IDIAP**

$$A(\tau) = \max_{d_i} \mathrm{FMR}_{d_i}(\tau) - \min_{d_j} \mathrm{FMR}_{d_j}(\tau)$$

**NIST**

$$A(\tau) = \frac{\max_{d_i} \mathrm{FMR}_{d_i}(\tau)}{\min_{d_j} \mathrm{FMR}_{d_j}(\tau)} \quad \forall d_i, d_j \in \mathcal{D}$$

**AWS**

$$A(\tau) = \sum_{d \in \mathcal{D}} \left| \log_{10} \frac{\mathrm{FMR}_d(\tau)}{\mathrm{FMR}^\dagger(\tau)} \right|$$

**IDEMIA**

$$A(\tau) = \frac{\max_{d_i} \mathrm{FMR}_{d_i}(\tau)}{\mathrm{FMR}^\dagger}$$

**MDTF (GINI)**

$$A(\tau) = \frac{\sum_i \sum_j |\mathrm{FMR}_{d_i}(\tau) - \mathrm{FMR}_{d_j}(\tau)|}{2n^2 \mathrm{FMR}^\diamond} \frac{n}{n-1}$$

# Two demographic summary measures

$$A(\tau) = \frac{\max_{d_i} \text{FMR}_{d_i}(\tau)}{\text{FMR}^{\dagger}}$$

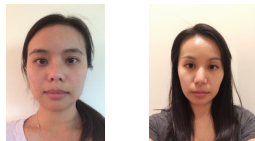Worst-case error rate over all demographic groups divided by the geometric mean

$$A(\tau) = \frac{\sum_i \sum_j |\text{FMR}_{d_i}(\tau) - \text{FMR}_{d_j}(\tau)|}{2n^2 \text{FMR}^{\diamond}} \frac{n}{n-1}$$

Mean absolute error rate difference over all demographic groups, divided by the arithmetic mean

$$x^{\dagger} = \left(\prod_i x_i\right)^{1/n}$$

$$x^{\diamond} = n^{-1} \sum_i x_i$$

# False Match Rates have bigger demographic variations

**Visa**  **Border**

| A: Lowest false match rates often in E. European men | B: Highest false match rates in older W. African women | C: False match rates 10-100 times higher | D: Economists standard measure "Gini" is much higher |

| Algorithm | Submission Date | FNMR Overall | FMR Min | FMR Max | FMR Max/GeoMean | FMR Gini |
|---|---|---|---|---|---|---|
| idemia_009 | 2022-07-27 | 0.0020 | 0.00027 C.America M (50-65] | 0.00641 W.Africa F (65-99] | 8.9[3] | 0.38[1] |
| cogent_007 | 2022-04-11 | 0.0034 | 0.00003 E.Europe M (35-50] | 0.00868 W.Africa F (65-99] | 25.6[187] | 0.61[108] |
| paravision_010 | 2022-02-02 | 0.0026 | 0.00000 S.Asia M (35-50] | 0.00219 W.Africa F (65-99] | 21.8[118] | 0.62[141] |
| s1_005 | 2022-06-17 | 0.0019 | 0.00002 E.Europe M (35-50] | 0.01039 W.Africa F (65-99] | 29.1[240] | 0.63[166] |
| cognitec_004 | 2022-02-10 | 0.0088 | 0.00005 E.Europe M (20-35] | 0.02211 W.Africa F (65-99] | 30.2[250] | 0.65[230] |
| sensetime_007 | 2022-06-17 | 0.0015 | 0.00004 E.Europe M (20-35] | 0.01565 W.Africa F (65-99] | 34.4[272] | 0.67[265] |
| rankone_013 | 2022-07-21 | 0.0021 | 0.00010 E.Europe F (12-20] | 0.03608 W.Africa F (65-99] | 52.1[325] | 0.76[334] |
| megvii_005 | 2022-03-28 | 0.0018 | 0.00001 E.Asia M (20-35] | 0.01059 W.Africa F (65-99] | 102.8[353] | 0.81[352] |

Source: https://pages.nist.gov/frvt/html/frvt11.html

**Conclusions:**
1. False negative rates vary greatly across demographic groups (age, gender, region-of-birth)
2. Some developers have improved

# Demographics: A False Positive Anecdote

NIST

**BLACK GIRL BANNED FROM MICHIGAN SKATING RINK BECAUSE FACIAL RECOGNITION SOFTWARE MISIDENTIFIED HER**

by Cedric 'BIG CED' Thornton ⏱ July 16, 2021 👁 4948



*(Image: Fox 2 Detroit)*

A young Black girl was kicked out of and banned from a skating rink in Michigan through no fault of her own. The girl was been banned due to facial recognition software that misidentified her as someone else.

https://www.zdnet.com/article/backlash-to-retail-use-of-facial-recognition-grows-after-michigan-teen-kicked-out-of-skating-rink-after-false-match/

**NIST**

THANKS    PATRICK.GROTHER@NIST.GOV    FRVT@NIST.GOV



NIST INTERAGENCY REPORT 8429
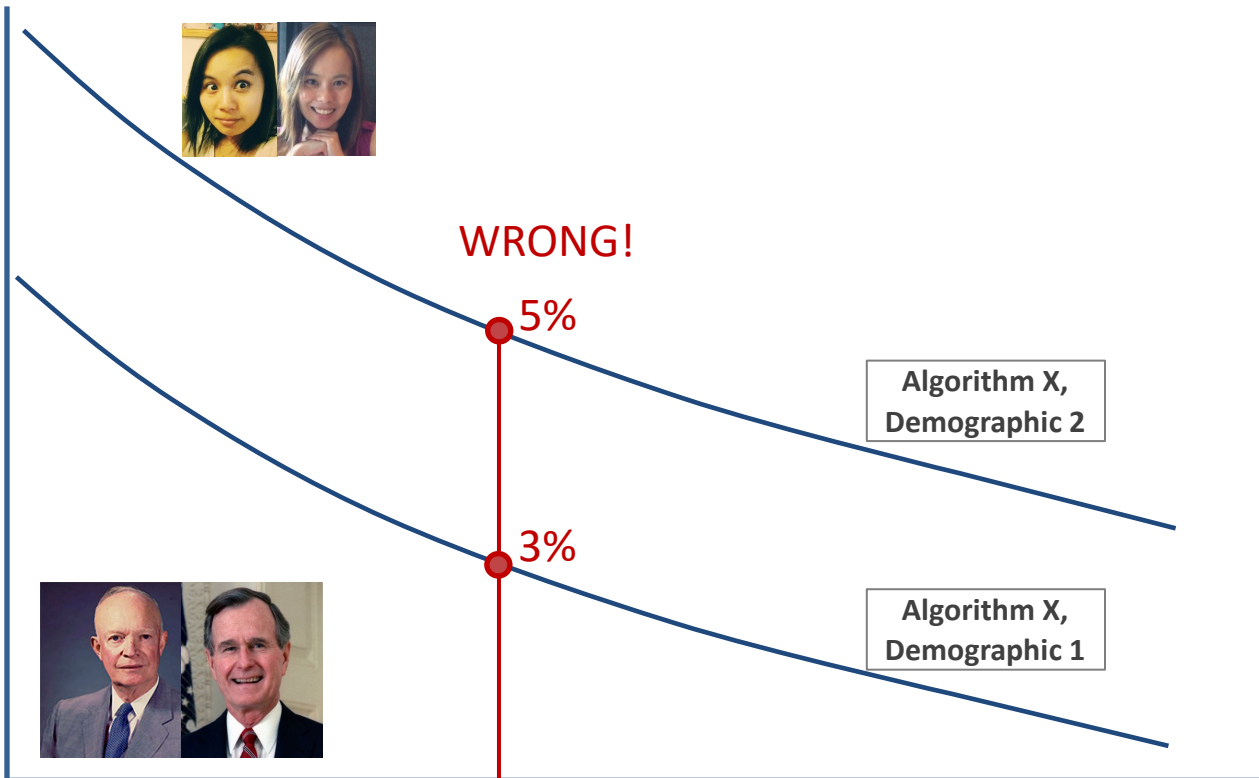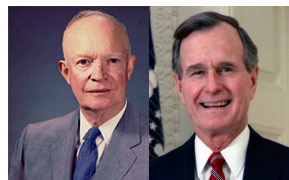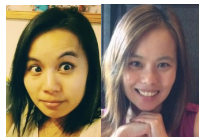SUMMARIZING DEMOGRAPHIC
DIFFERENTIALS



ISO/IEC 9868 WD7
PASSIVELY CAPTURED SUBJECTS



ISO/IEC 19795-10 WD4
QUANTIFYING DEMOGRAPHIC EFFECTS

**FNMR**
**False non-match rate**

Proportion of genuine comparisons producing score below threshold, T.

See ISO/IEC 19795-1

Δ FMR

Δ FNMR

T = 3.142

T = 3.142

Algorithm X, Demographic 2

Algorithm X, Demographic 1

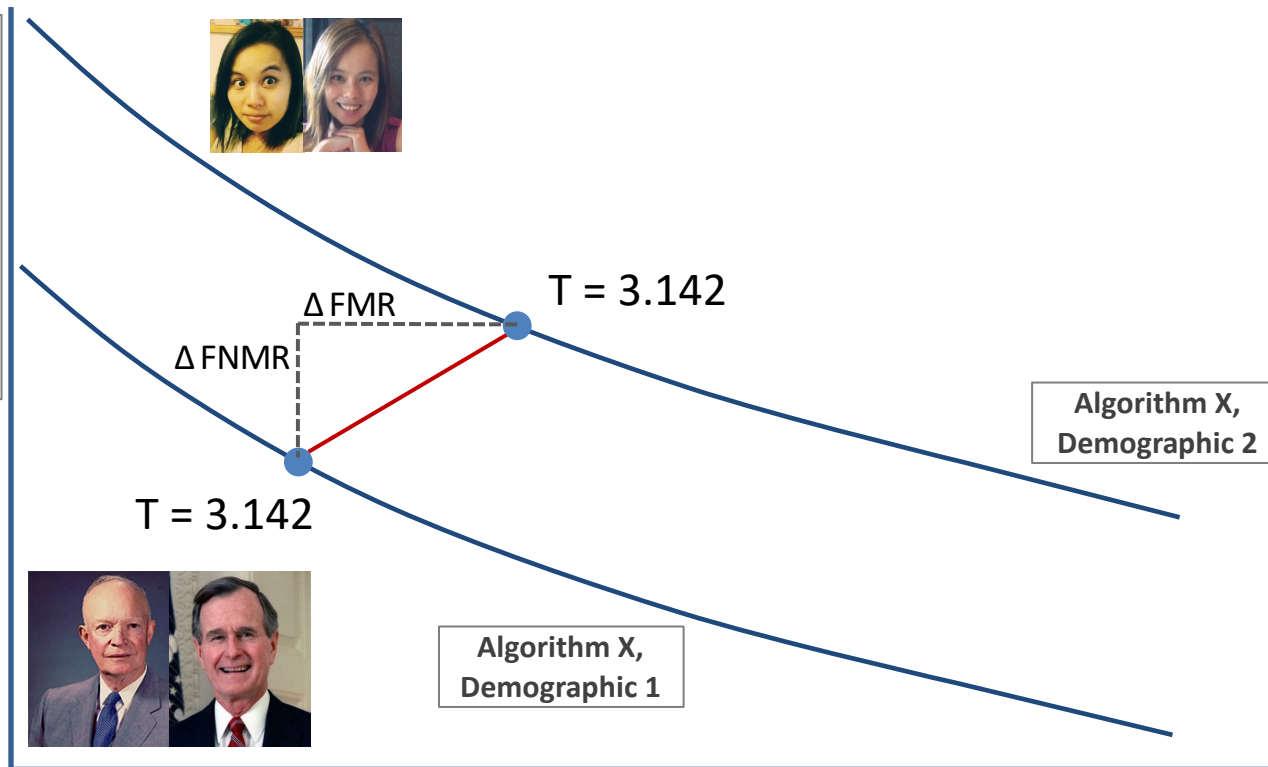Low FMR values achieved with higher, i.e. more stringent, thresholds.

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.

29