

NISTIR 7679

NIST Special Database 32
Multiple Encounter Dataset I (MEDS-I)

Data Description Document

Steven Curry
Drew Founds
Joe Marques
Nick Orland
MITRE Corporation

Craig Watson
Information Technology Laboratory

NISTIR 7679

NIST Special Database 32
Multiple Encounter Dataset I (MEDS-I)

Data Description Document

Steven Curry
Drew Founds
Joe Marques
Nick Orland
MITRE Corporation

Craig Watson
Information Technology Laboratory
Information Access Division

December 2009



U.S. Department of Commerce

National Institute of Standards and Technology
Patrick D. Gallagher, Deputy Director

Acknowledgement

This dataset is being released (as prepared by MITRE Corporation) to support the NIST Multiple-Biometric Evaluation 2010 (MBE). In addition, this dataset is available to any user interested in biometric research. The sponsor of this joint effort and provider of the data is the Federal Bureau of Investigation (FBI).

Disclaimer

Specific hardware and software products identified in this report were used in order to process and analyze the dataset described in this document. In no case does identification of any commercial product, trade name, or vendor, imply recommendation or endorsement by the National Institute of Standards and Technology, nor does it imply that the products and equipment identified are necessarily the best available for the purpose.

MITRE

Multiple Encounter Dataset (Deceased Persons)

Data Description Document

Steven Curry

Drew Founds

Joe Marques

Nick Orland

Version 3.0

November 2009



MEDS: Multiple Encounter Dataset (Deceased Persons)

Data Description Document

Steven Curry
Drew Founds
Joe Marques
Nick Orlans

Version 3.0
November 2009

Sponsor: FBI CJIS
Dept. No.: G123
Contract No.: J-FBI-07-164
Project No.: 1408FC09
Downgrade UNCLASSIFIED
Derived By:
Declassify On:

The views, opinions and/or findings contained in this report are those of The MITRE Corporation and should not be construed as an official government position, policy, or decision, unless designated by other documentation.

This technical data was produced for the U. S. Government under contract J-FBI-07-164 and is subject to the Rights in Data-General Clause 52.227-14 (JUNE 1987).

©2009 The MITRE Corporation.
All Rights Reserved.

1 Overview

Multiple Encounter Dataset (MEDS) is a test corpus organized from an extract of submissions of deceased persons with prior multiple encounters. MEDS is provided to assist the FBI and partner organizations refine tools, techniques, and procedures for face recognition as it supports Next Generation Identification (NGI), forensic comparison, training, and analysis, and face image conformance and inter-agency exchange standards. The MITRE Corporation (MITRE) prepared MEDS in the FBI Data Analysis Support Laboratory (DASL) with support from the FBI Biometric Center of Excellence.

This paper describes basic characteristics of the face images and some relevant image quality characteristics that pertain to both collection practices and the calibration and evaluation of face recognition technology.

All 685 original submissions contain at least one Type 10 record, the record type within the ANSI/NIST-ITL 1-2007 file format reserved for face images and Scars, Marks, and Tattoo (SMT) images. ¹ Two submissions lacked usable face images and SMT images were removed from the dataset, leaving 683 submissions in the collection. As described in table 1-2, 686 images are frontal or “near frontal.” Near frontal images are within 10 degrees horizontal of full frontal as estimated by visual inspection. No estimate for pitch, roll, and yaw (3D pose angles) has been determined although some variations exist. Profiles or “near profiles” are likewise within 10 degrees of full profile, although “two-eyed” profiles may be closer to 20 degrees-25 degrees off full profile, and there is at least one quarter profile (approximately 45 degree). Tables 1-1 and table 1-2 below summarize the submissions in terms of the number of encounters, subjects, and Type 10 records.

Table 1-1: Number of Encounters, Subjects, and Submissions

Number of Encounters	Subject Count	Number of Files (submissions)
1	248	248
2	70	140
3	25	75
4	13	52
5	7	35
6	2	12
7	6	42

¹ American National Standard for Information Systems – Data Format for the Interchange of Fingerprint, Facial, and other Biometric Information – Part 1. NIST Special Publication 500-271, May, 2007. Online: <http://fingerprint.nist.gov/standard/Approved-Std-20070427.pdf>

Number of Encounters	Subject Count	Number of Files (submissions)
8	4	32
9	3	27
10	2	20
	380	683

Table 1-2: Number of Type 10 Image Files (Face)

Type 10 Records (per submission)	Count	Comments
1	683	683 frontal or near frontal faces (after removal of one “NO PHOTO” and one other)
2	26	24 profiles, two “two-eyed” profiles
3	3	One frontal, one profile, one quarter profile (one SMT removed)
Type 10 Face Images	712	

The final number of face images in the collection is 712. The number of submissions is 683, and the number of subjects is 380.

2 Methodology

MITRE parsed and examined the submission files using a combination of government, commercial, and custom Electronic Biometric Transmission Specification (EBTS) parsing and reporting tools to help verify consistent results. The following tools were used:

- Universal Latent Workstation, v 5.4.0 (Government Off the Shelf (GOTS)/Noblis ²), manual inspection and viewing
- Sanitizer (GOTS/MITRE), batch processing and sanitization
- EBTSExtract (GOTS/MITRE), batch extraction and reporting
- ParseEFTS (COTS/Aware ³), batch extraction, verification

² <http://www.noblis.org>

³ <http://www.aware.com>

Picasa ⁴, PittPatt ⁵, OpenCV ⁶ (COTS and opensource), reference tools for face detection
Microsoft Excel (COTS), data tabulation

3 Description of Data

3.1 Image Sizes

The sizes of the images and resolution of the subject vary significantly. The range of image sizes and their occurrences are described in table 3-1.

Table 3-1: Type 10 Image File Sizes

Type Image Size (pixel width x height)	Count
240 x 240	75
300 x 300	2
300 x 400	1
384 x 480	156
476 x 596	2
480 x 600	365
486 x 621	1
494 x 620	1
618 x 794	1
640 x 480	1
657 x 389	1
736 x 858	1
768 x 960	93
793 x 981	1
960 x 1280	2
1280 x 960	6
1824 x 1170	1
2272 x 1704	2
Total Face Images	712

3.2 Face Resolution

Consistent resolution of the faces requires consistent framing of the subject for a particular image size. The framing of the subject in the MEDS images varies and

⁴ <http://picasa.google.com>

⁵ <http://pittpatt.com>

⁶ <http://sourceforge.net/projects/opencvlibrary/>

⁷ This is after 90 degree rotation, formerly these two images were 1280 x 960

in some instances the full face is not visible. For the frontal and near frontal images, MITRE estimated interocular distances based on the outputs from automated face detection using PittPatt⁸. The results presented in Figures 3-1 and 3-2 are based on 672 “both eye” detection events from the 685 frontal or near frontal images. The results are presented based on the automated outputs and have not been manually verified. Table 3-2 illustrates the reported misses. Note these are not necessarily errors as some images are profiles.

Table 3-2: PittPatt Misses

Misses	Count
Left Eye Misses	34
Right Eye Misses	20
Both Eyes Missed	13

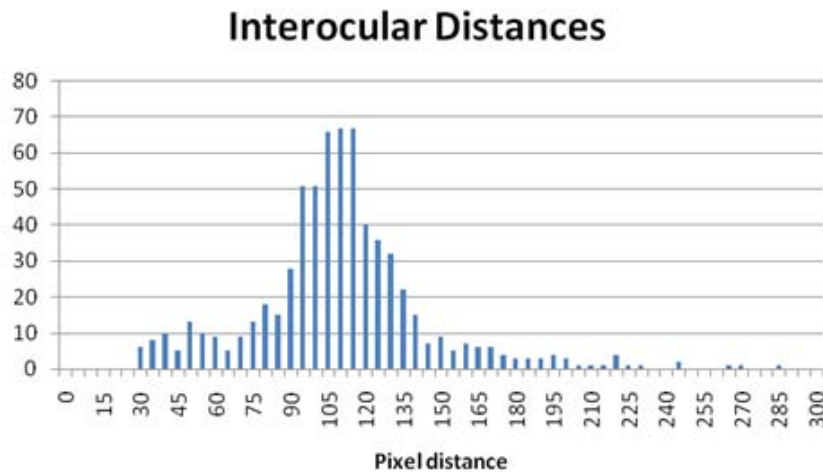


Figure 3-1: Estimated Interocular Distances (pixel distance)

⁸ Pittsburgh Pattern Recognition Face Tracking and Detection Software Development Kit, version 2.2. Online reference <http://pittpatt.com>.

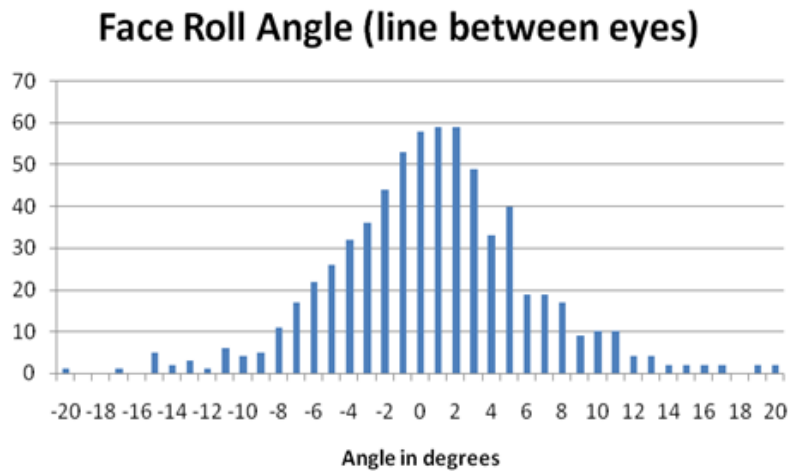


Figure 3-2: Face Roll Angle (degrees from horizontal)

The Information Technology Laboratory (ITL) specification has no image resolution requirements that are applicable to Type 10 records. *“Facial/mugshot, SMT, and iris images rely on the total number of pixels scanned and transmitted and are not dependent on the specific scanning resolution used.”* (Section 6, page 11).

Subject Acquisition Profiles (SAPs) provide recommended ranges and *desired* resolutions; is this sufficient for NGI and derived requirements for automated binning, workflow optimization, and human examination?

3.3 Age and Longitudinal Data

The MEDS data contains subjects with diverse ages and provides examples of repeated observations of the same individual over time. The ages over all submissions are summarized below in figure 3-4. The time between first and last encounters are summarized in figure 3-5, and reported in full in the accompanying text file “medsDurations.txt”. The 248 durations of zero are single encounters. The 0 to 0.1 interval represents durations of less than 30 days.

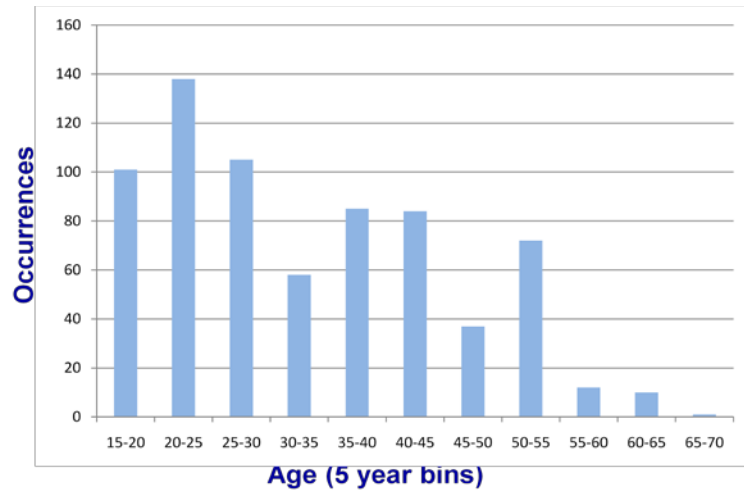


Figure 3-4 Age Ranges of Subjects (for each of the 683 submissions)

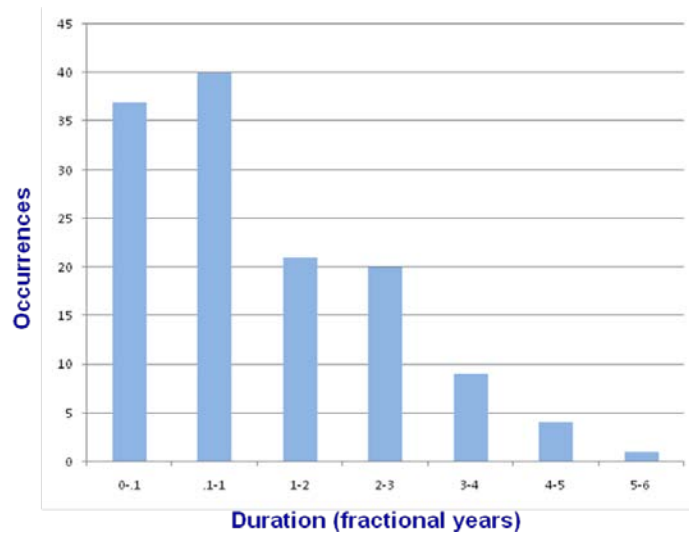


Figure 3-5 Duration (between first and last encounter, fractional years)

3.4 User-Defined and Type 10 Textual Data

The names of the subjects and the circumstances of their encounter or arrest are not relevant to biometric performance questions and hence are excluded for privacy reasons. However, some user-defined fields pertaining to acquisition standards and the subject's physical appearance (that are also evident from the photos) are maintained as potentially useful metadata. The fields maintained are listed below. The values are provided in the associated file 'medsMetaData.csv'.

Table 3-2: MEDS Metadata

Field	Field name	Comments
2.022	Date of Birth	User-defined
2.024	Sex	User-defined
2.025	Race	User-defined, NCIC codes
2.029	Weight	User-defined
2.031	Eye Color	User-defined, NCIC codes
2.032	Hair Color	User-defined, NCIC codes
10.003	Image Type	124 of the 712 faces are incorrectly marked "RGB"
10.005	Photo Date	124 of the 712 faces are missing photo date, for these images the date of arrest (2.045) is used instead. The missing photo dates occurs on the same 124 images that incorrectly used color space for image type.
10.006	Horizontal Line Length	
10.007	Vertical Line Length	
10.025	Subject Pose Angle	

4 Notes

The following are miscellaneous notes and comments on the MEDS images.

4.1 Missing Camera Information and Lack of EXIF

Subject acquisition source (10.023) was not used on any of these data, and no Exchangeable Image File Format (EXIF) data was available for any images. There is a file format conflict that currently prevents the use of EXIF header information in JPEG File Interchange Format (JFIF) containers as specified within the ITL 1-2007 standard. ITL requires JPEG encoded data be wrapped in JFIF containers, which are marked by Tag APP0. The JFIF tag appears after the start of the image (SOI) marker, however, the EXIF specification states that the EXIF tag APP1 appears after the SOI marker. Thus if EXIF data exists with jpeg images, current generation ITL parsers will bypass it to be strictly ITL conformant.

This topic warrants additional attention from practitioners and the standards community on how to preserve and leverage EXIF information, providing more reliable device level information for image quality (e.g., device, color spaces, and compression).

4.2 Rotations

The following two images were rotated:

S184-01_t10_1 was rotated 90 CCW to be upright

S184-01_t10_2 was rotated 90 CW to be upright

4.3 Names in Images

Four images contained the subject's name in the image. These images were manually modified to obscure the name. Obscuration was done using Adobe Photoshop's "clone" brush by manually brushing adjacent colors over the text.

The following four images were modified in this manner:

S018-01-t10_01 (name rendered in image)

S131-01-t10_01 (obfuscated name on mechanics shirt)

S180-01-t10_01 and S180-01-t10_02 (obfuscated name tag)

4.4 Non-Face Images

Three images were removed as non-faces or out of scope as described below:

One image was marked as 'FACE" and but contained a blank image "NO PHOTO"

One image was an arm (SMT) image

One image was unusable

4.5 Scars, Dirt, Bruises, Bandages, and Other Factors.

Several examples of scarring and bruising are evident. Some of the subjects have bandages, presumably from recent injury, while others have dirt or grime that may affect skin tone and or texture. These factors are not fully documented here and in some cases may be difficult to interpret due to the lack of ground truth for comparison with a "clean" sample.

5 Baseline Face Detection

MITRE conducted a baseline reference for face detection performance using Google Picasa (beta 3). The Picasa client is a photo organizer available for no-cost download and is designed to work in conjunction with online storage capabilities provided by Google. The experiment was conducted on an isolated network and involved no contact or use of server based functions.

The 712 face images were loaded into Picasa and the “show photos with faces” face detection filter was applied. For all 712 face images, there were 13 missed detections of frontal or near frontal faces and 27 missed detections of all profile and near profile faces. This represents an overall detection rate of 94.3 %.

Eliminating the profile and near profile faces, 685 frontal or near frontal faces generate 13 missed detections with a 98.1 % detection rate. The 13 misses are:

Table 5-1: Picasa Missed Detections

Number	Image Filename
1	S013-01_t10_1
2	S073-01-t10_2
3	S119-02_t10_1
4	S182-01_t10_1
5	S184-01_t10_1
6	S184-02_t10_1
7	S187-02_t10_1
8	S202-01_t10_1
9	S230-01_t10_1
10	S238-02_t10_1
11	S247-01_t10_1
12	S264-01_t10_1
13	S340-01_t10_1