

Biometrics Testing and Statistics

Introduction

Determining the best biometric system for a specific operational environment and how to set up that system for optimal performance requires an understanding of the evaluation methodologies and statistics used in the biometrics community. This document provides a baseline testing and statistics review, thus enabling appropriate analysis of available research reports. This document is intended to further the understanding of a general audience and is not intended to replace or compete with sources that may be more technically descriptive/prescriptive such as those under development by standards bodies such as INCITS and ISO/IEC. Detailed information on how to properly perform performance evaluations is beyond the scope of this document.

Evaluation Types

Performance evaluations of biometric identification technology are divided into three overlapping categories with increasing complexity in uncontrolled variables: technology, scenario, and operational.¹ A thorough evaluation of a system for a specific purpose starts with a Technology Evaluation, followed by a Scenario Evaluation, and finally an Operational Evaluation.

The primary goal of Technology Evaluations is to measure the performance of biometric systems, typically only the recognition algorithm component. They are repeatable and usually short in duration. Technology Evaluations are usually performed using standard datasets collected previous to testing. In general, results from a Technology Evaluation show specific areas that require future research and development (R&D) and provide performance data that is useful when selecting algorithms for scenario evaluations. An example of a Technology Evaluation is the Face Recognition Vendor Test.²

The primary aim of Scenario Evaluations is to measure performance of a biometric system operating in a particular application. For example, testing biometrics for access control purposes at a mock doorway in a laboratory. Each tested system normally would have its own acquisition sensor and would thus receive and produce slightly different data. For this and other reasons, Scenario Evaluations are not always completely

National Science and Technology Council (NSTC)

Committee on Technology

Committee on Homeland and National Security

Subcommittee on Biometrics



repeatable. Scenario Evaluations usually take a few weeks to complete because multiple trials (and for some Scenario Evaluations, multiple trials of multiple subjects/areas) must be completed to ensure adequate habituation of the end users (if the scenario calls for it) and to achieve a statistically relevant number of samples. Results from a typical Scenario Evaluation show areas that require additional system integration and provide performance data on systems for the application tested. An example of a Scenario Evaluation is the UK Biometric Product Testing.³

At first glance, an Operational Evaluation appears very similar to a Scenario Evaluation, except that the test is conducted at the actual site using actual end users, a subset of the end users, or a representative set of subjects. Rather than testing for performance (which is difficult, if not impossible, to do in some operational evaluations), Operational Evaluations typically aim to determine the workflow impact caused by the addition of a biometric system. Operational Evaluations are typically not repeatable. Operational Evaluations can last from several weeks to several months because the evaluation team must first examine workflow performance prior use of the technology and again after users are familiar with the technology. An accurate analysis of the benefit of the new technology requires a comparison of the workflow performance before and after use of the technology.

In an ideal three-step evaluation process, Technology Evaluations are first performed on all applicable technologies that could conceivably meet requirements. The technical community then uses the results to plan future R&D activities, while potential users use the results to select promising systems for application-specific Scenario Evaluations. Results from the Scenario Evaluation(s) will enable users to determine the best system for their specific application and to have a good understanding of how it will operate at the proposed location. This performance data, combined with workflow impact data from subsequent Operational Evaluations, will enable decision makers to develop a solid business case for potential installations.

So for those analyzing evaluation reports, it is important to determine which type of evaluation occurred and its relevance to an intended application. Generally, technology evaluation reports contain information relevant to most intended applications of a given biometric, while operational evaluation reports are generally only useful if the intended application is very closely related to what was tested.



Biometric Evaluation Terms

Biometric terms such as recognition, verification and identification are sometimes used interchangeably. This is not only confusing but incorrect as each term has a different meaning.

- Verification occurs when the biometric system attempts to confirm an individual's claimed identity by comparing a submitted sample to one or more previously enrolled templates.
- Identification occurs when the biometric system attempts to determine the identity of an individual. A biometric is collected and compared to all the templates in a database. Identification is "closed-set" if the person is assumed to exist in the database. In "open-set" identification, the person is not guaranteed to exist in the database. The system must determine if the person is in the database. A "watchlist" task is an example of "open-set" identification.
- Recognition is a generic term and does not necessarily imply either verification or identification. All biometric systems perform "recognition" to "again know" a person who has been previously enrolled.

This section provides in-depth, clearly defined descriptions of these tasks. To help explain them, a hypothetical face recognition system must be introduced. This hypothetical face recognition system can compare one image to another and provide scores (*similarity scores*^a) for each comparison. For our example system, the similarity scores range from 0.0 to 1.0, with a 1.0 score being an exact match. The system also has a user-set "threshold" that the system uses to make a matching decision. Although the examples in this section use face recognition, the tasks and associated performance measures are the same as for other biometric types.

^a Not all biometric systems use similarity scores for comparisons. Some use difference scores, hamming distances, etc. For the purposes of this non-technical paper, the basic concept is essentially the same - mathematically comparing two biometric templates in order to make a matching decision.



Verification

In the verification task, an end user must first make a claim as to his/her identity (e.g., I am John Q. Public) and the biometric system then determines if the end-user's identity claim is true or false. A good example is verifying an end user's identity, frequently represented by a username, by requiring a password prior to providing access to his/her account on a computer system. Figure 1 gives a visual example where the gentleman on the right makes a claim that he is the gentleman on the left. For this example, assume these are pictures of the same individual.

**CORRECT
VERIFICATION
CLAIM**



Figure 1: Correct Verification Claim.

Assume that the example face recognition system produces a similarity score of 0.93 for this verification trial. (Remember that our demonstration face recognition system works on a 0.0 to 1.0 scale with 1.0 being an exact match.) Also assume that the system's verification threshold was set at 0.90. Since 0.93 is higher than 0.90, the system in this example has correctly determined that the gentleman in the right picture is the same as the gentleman in the left picture. This is called a true accept or correct verification.

Now assume that the same individual in Figure 1 makes the same claim, except this time the system's verification threshold is set at 0.95. In this case, the demonstration face recognition system will not make a correct decision.^b

If we run many trials with this gentleman, as well as other correct matches, we will know the rate^c at which legitimate end users are correctly verified by the system. This is called the true accept or correct verification rate.

^b This situation is referred to as a *false reject*.

^c Technically, these tests will produce a statistical estimate of the actual rate. For simplicity sake, the term "rate" is used in this introductory document.



Figure 2 shows a different verification claim. In this example, the gentleman on the right claims to be the gentleman on the left. Obviously, this is not the case. Assume that the system returns a similarity score of 0.86. Let us also assume that the system's verification threshold was set at 0.9. In this example, the face recognition system determines that the gentleman on the right is not the gentleman on the left.

**FALSE
VERIFICATION
CLAIM**



Figure 2: False Verification Claim.

Now let us look at the case where the same individual in Figure 2 makes the same claim, but the system's verification threshold is set at 0.85. In this case, the system incorrectly verifies that the gentleman is the gentleman in the system. This error is called a false accept. If many trials are run with incorrect claims, the rate at which the system incorrectly matches an imposter individual to another individual's existing biometric will be known. This is called the false accept rate.

Ideally, biometric systems would always provide a probability of verification of 100% with a false accept rate of 0%.^d Unfortunately, that is not possible; so system administrators must compromise by setting the system's threshold at an optimum value for their given application. Determining the threshold can be difficult because the verification rate and false accept rate are not independent variables.^e If the threshold in the example face recognition system is raised, the verification rate decreases, but the false accept rate also decreases. If the threshold in the example system is lowered, the verification rate rate increases, but the false accept rate also increases. Plotting verification

^d From a statistical standpoint, neither of these results is even possible. Someone may run a test with no observed errors, but they statistically wouldn't have a 100% reliable system. All documented results should meet basic statistic principles.

^e This relationship is similar to that of a metal detector. By adjusting the threshold, security personnel increase the chances of it alarming on larger or smaller metal items.



accept rates against the associated false accept rates, called a Receiver Operating Characteristic (ROC) curve, allows for a visualization of this trade-off relationship. Figure 3 is a sample of a verification ROC (with fabricated numbers for example purposes). Varying the system's threshold moves the operating point along its ROC curve.

Receiver Operating Characteristic

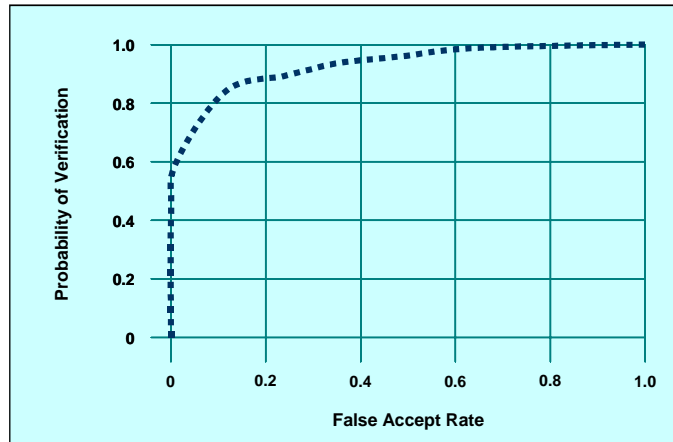


Figure 3: Example Verification ROC.

Open-Set Identification

In open-set identification (sometime referred to as a watchlist application), the biometric system determines if the individual's biometric template matches a biometric template of someone in the database. The individual does not make an identity claim and, in cases of covert identification, does not personally interact with the system whatsoever. Examples of this task might be comparing biometrics of visitors to a building against a terrorist database, or comparing a biometric of a "John Doe" in a hospital to a missing person's database. Figure 4 shows an image of a gentleman as an input to the example face recognition system.

The system first compares the submitted image to each image in the database. Assume that the similarity score for each comparison is 0.6, 0.86, 0.9, and 0.4 (respectively). Also assume that the system's watchlist threshold is set at 0.85. In this example, the face recognition system sounds an alarm each time one or more of the similarity scores is higher than the threshold. Since an alarm sounded, the system user would look more closely at the similarity scores to see which image attained the highest score, which would be the system's best guess at the identity of



the subject in the input image. We can easily see that it is correct. (This description is only concerned with the top match, so the fact that a second comparison also had a similarity score higher than the threshold is irrelevant.) This example produced what is called a *correct detect and identify*.

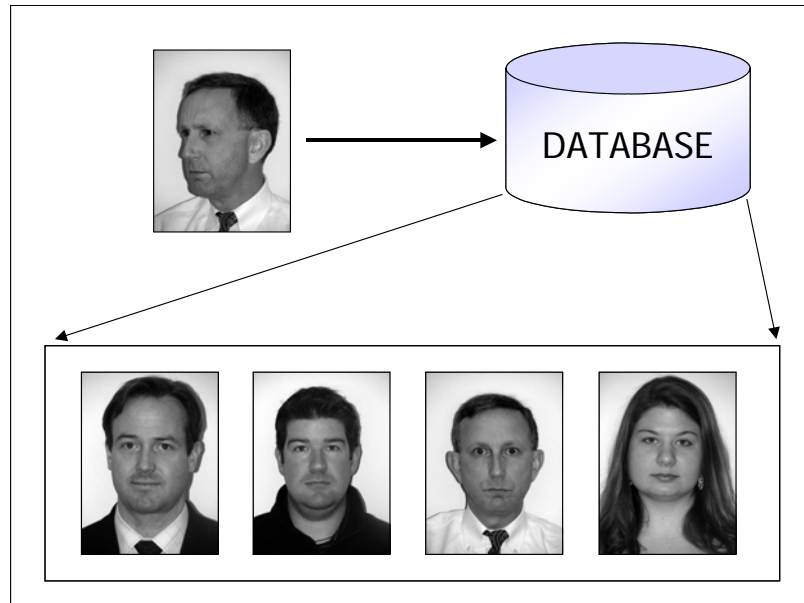


Figure 1: Watchlist Example 1.

Consider again the example shown in Figure 4, except this time the watchlist threshold is 0.95. In this case, the face recognition system does not sound an alarm because none of the similarity scores (0.6, 0.86, 0.9, and 0.4) are above the system's threshold. Since there was no alarm, there would be no reason to look further at the similarity scores. Thus, for this example, the demonstration face recognition system did NOT produce a *correct detect and identify*.

Taking a final look at the example shown in Figure 4, assume that the similarity score for each comparison is 0.6, 0.86, 0.8, and 0.4, respectively, and the watchlist threshold is 0.75. In this example, the system sounds an alarm as one or more of the similarity scores are higher than the threshold. The system user would look more closely at the similarity scores and see that the second individual has the highest score. In this example, an alarm correctly sounded (as the subject is in the database), but the demonstration face recognition system did not correctly choose the identity of the gentleman as the top-ranked match. Thus, the system did NOT produce a *correct detect and identify*.



If we run many trials, we will know how often the system will return a correct result. A correct result occurs when an individual who is in a database causes a system alarm AND is properly identified in an open-set identification (watchlist) application. This is called the *Detect and Identification Rate*.

Now consider an alternative setup where the input does not have a corresponding match in the database, as shown in Figure 5. In this example, an image of a lady is the input to our example face recognition system, which must determine if this individual is in the database.

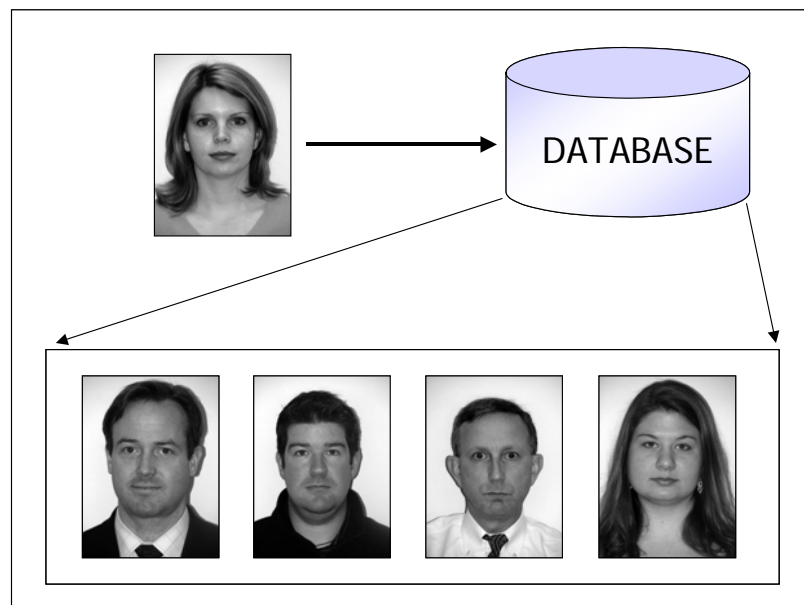


Figure 2: Watchlist Example 2.

The system first compares the input image to each image in the database. Assume that the similarity score for each comparison is 0.7, 0.8, 0.4, and 0.6, respectively, and the system's watchlist threshold is set at 0.85. In this example, an alarm will not sound, as none of the similarity scores are higher than the threshold.

Now consider the same example with a threshold set at 0.75. In this case, an alarm sounds because one of the similarity scores is higher than the threshold. This is an incorrect alarm, because the lady in the input image is not in the database. This is called a *false alarm*. If we run many trials with subjects who are not in the database, we will know how often the system will return an incorrect alarm, i.e., the *false alarm rate*.



Because biometric systems cannot provide a detection and identification rate of 100% with a false alarm rate of 0%, system administrators must set the system's threshold at an optimum value for the given application and the tradeoffs of correctly identifying subjects versus false alarms. If the watchlist threshold in the example system is raised, the identification rate decreases, but the false alarm rate also decreases. If the watchlist threshold is lowered, the identification rate increases, but the false alarm rate increases. Plotting the identification rates and the associated false alarm rates, also called a *Receiver Operating Characteristic (ROC)*, allows for a visualization of this trade-off relationship. These are sometimes referred to as a Watchlist ROC or an Identification ROC to help differentiate it from a verification ROC. Figure 6 is an example watchlist ROC (with fabricated numbers).

Watchlist ROC

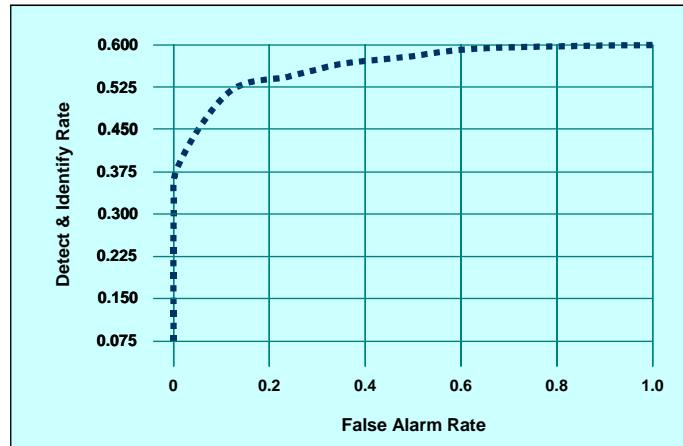


Figure 1: Example Watchlist ROC.

Database size is important to watchlist performance. The Face Recognition Vendor Test (FRVT) 2002⁴ showed that watchlist performance for face recognition systems decreases as the size of the database increases. (Effectively, the curve in Figure 6 will lower as the size of the database increases.) When quoting open-set identification performance, it is important to also state the database size.

In practice, the open-set identification task is much more difficult for biometric systems (and presumably for human operators) than the verification task. When discussing a specific application, it is critical to think in terms of the proper task and the associated



statistics. Failure to do so will lead to significant confusion and errors.

Closed-Set Identification

Closed-set identification is where every input image has a corresponding match in the database. In practice, there are very few applications that operate under the closed-set identification task. Even the FBI's Integrated Automated Fingerprint Identification System (IAFIS) operates as a watchlist -- an open-set identification task. However, these statistics are routinely found in research and evaluation reports, as they are a good measure of showing general strengths and weaknesses.

In the closed-set identification task, a biometric template of an individual is presented to the biometric system, as shown in Figure 4. Again, it is known that the person is in the database. The example face recognition system first compares the input image to each image in the database. Let us assume that the similarity score for each comparison is 0.6, 0.86, 0.9, and 0.4, respectively. In this example, the correct match has the top similarity score. If we run the same trial for all subjects in the database, we will know how often the system will return a correct result with the top match, which is termed the identification rate at rank 1.

Still referring to the example shown in Figure 4, assume that the similarity score for each comparison is 0.6, 0.4, 0.8, and 0.86, respectively. In this case, the correct match is the second highest similarity score. If we run the same trial for all subjects in the database, we will know how often the system will return a correct result in either the top or second ranked score. (We do not necessarily care if they are in the top or second rank specifically, just that they are in one of those positions.) This is termed the identification rate at rank 2.

These two examples show a trend for how to show identification performance graphically. The probability of correct identification at rank 20 means, what is the probability that the correct match is somewhere in the top 20 similarity scores? A Cumulative Match Characteristic (CMC) curve shows the probability of identification for numerous ranks. Figure 7 is an example CMC (with fabricated numbers for example purposes).



Cumulative Match Characteristic

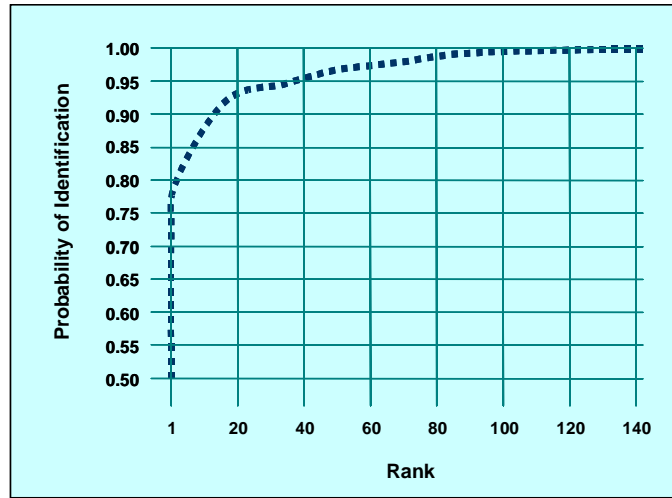


Figure 2: Example Cumulative Match Characteristic Curve.

One key feature of a CMC is that, in a plot that includes all possible ranks (e.g., if the database has 140 people, and the CMC goes through rank 140), the probability of identification is 100% at the highest (140 in this example) rank. This is true because every input is in the database (otherwise this is an open-set identification task instead of a closed-set identification task), and it is showing the identification rate for the entire database.

Just as in the watchlist task, it is important to state the size of the database when describing a CMC curve. The probability of correct identification at rank 10 for a 100-person database would be much better than the probability of correct identification at rank 10 for a 10,000-person database (all other factors being the same).

Failure to Acquire

This document has described the three biometric tasks and their associated performance measures. However, there is another measure that may also be of interest because it affects all three biometric tasks. The Failure to Acquire rate is the rate at which a biometric system fails to capture and/or extract information from an observation. Numerous issues, including device/software malfunction, environmental concerns, and human anomalies (e.g., amputees not able to use hand geometry system, bricklayers with worn fingerprints, etc.), can cause a Failure to Acquire. For some biometric systems, or for certain applications, the Failure to Acquire rate could be quite high.



Different evaluations deal with this issue in different ways. Some (as in the examples above) force systems to produce similarity scores, even if there was a Failure to Acquire. This, of course, produces lower performance measures. Others only show performance (usually referred to as False Match Rates^a and False Non-Match Rates^b) on properly acquired signatures and show the Failure to Acquire rate separately. This, of course, raises the performance measures. Neither approach is wrong; evaluators simply choose the method that shows performance according to how the system will be used operationally. When reviewing others' evaluations, potential users will need to determine which approach was applied.

Other Performance Statistics

Other statistics are sometimes used to show performance of biometric systems. These, listed below, are defined in the accompanying Glossary.

- Crossover Error Rate (CER)
- Detection Error Trade-off (DET)
- Difference Score
- Equal Error Rate (EER)
- Failure to Enroll (FTE)
- False Match Rate
- False Non-Match Rate
- Hamming Distance
- Throughput Rate
- True Accept Rate
- True Reject Rate
- Type I Error
- Type II Error

Other Types of Testing

Not all biometric tests are accuracy-based. A summary of the more common of these tests is described below.

^a The False Match Rate is equivalent to the False Acceptance Rate described in this paper.

^b The False Non-Match Rate is similar to the False Reject Rate (FRR) described in this paper, except the FRR includes the Failure to Acquire error rate and the False Non-Match Rate does not.



Acceptance Testing: "The process of determining whether an implementation satisfies acceptance criteria and enables the user to determine whether or not to accept the implementation. This includes the planning and execution of several kinds of tests (e.g., functionality, quality, and speed performance testing) that demonstrate that the implementation satisfies the user requirements."⁵

Conformity: "Fulfillment by a product, process or service of specified requirements"⁶

Conformity Evaluation: "Systematic examination of the extent to which a product, process or service fulfils specified requirements"⁶

Conformance Testing (or Conformity Testing): "Conformity evaluation by means of testing"⁶

Interoperability Testing: "The testing of one implementation (product, system) with another to establish that they can work together properly"⁷

Performance Testing: "Measures the performance characteristics of an Implementation Under Test (IUT) such as its throughput, responsiveness, etc., under various conditions"⁵

Robustness Testing: "The process of determining how well an implementation processes data which contains errors"⁵

Standards Activities

There are multi-part voluntary consensus standards for Biometric Performance Testing and Reporting under development by INCITS M1 and ISO/IEC JTC 1/SC 37. The first three parts of the INCITS American National Standard were approved by ANSI on October 25, 2005. These parts are:

INCITS 409.1-2005, American National Standard for Information Technology - Biometric Performance Testing and Reporting - Part 1: Principles and Framework. This multipart standard develops a common set of methodologies and procedures to be followed for conducting technical performance testing and evaluations. Included are guidelines that address issues regarding required test sizes, performance statistics, error reporting, and presentation of performance results. These procedures



can be incorporated in an "end-to-end" system approach or from an individual technical component perspective.

INCITS 409.2-2005, American National Standard for Information Technology - Biometric Performance Testing and Reporting - Part 2: Technology Testing and Reporting. This standard specifies procedures for conducting offline tests of the performance of biometric technologies.

INCITS 409.3-2005, American National Standard for Information Technology - Biometric Performance Testing and Reporting - Part 3: Scenario Testing and Reporting. This standard specifies requirements for scenario-based biometric testing and reporting.

A similar standard is under development at the international level, ISO/IEC FDIS 19795-1:2005. Part 1: Principles and Framework is up for ballot. 19795-1 has been developed from the UK Biometrics Working Group's Best Practices in Testing and Reporting Performance of Biometric Devices. The UK document was developed from two NIST primary source documents developed by NIST, a variety of evaluation reports and input from the Biometric Consortium's Working Group on Interoperability, Performance and Assurance.

Important Items to Keep in Mind

There are two key items to keep in mind while reviewing biometric performance evaluation reports. First, not all evaluation results are relevant. If an evaluation report, particularly for a Scenario or Operational Evaluation, does not match the user's intended application, the usefulness of the results will be significantly diminished. Second, biometric evaluation results have a very limited shelf life. Researchers continue to make significant progress in improving the performance of a biometric system so if the report is more than 9-18 months old, the results should not be considered conclusive, but merely used as a general guide and reference.

Document References

¹ P. Philips, A. Martin, C. L. Wilson, and M. Przybocki, "An Introduction to Evaluating Biometric Systems" (2000). <http://www.frvt.org/DLs/FERET7.pdf>.



² "Face Recognition Vendor Test," FRVT.org
<<http://www.frvt.org>>.

³ Tony Mansfield, Gavin Kelly, David Chandler, and Jan Kane, "Biometric Product Testing Final Report" 19 March 2001, CESG/BWG Biometric Test Programme <http://www.cesg.gov.uk>.
<<http://www.cesg.gov.uk/site/ast/biometrics/media/BiometricTestReportpt1.pdf>>.

⁴ P.J. Phillips, P. Grother, R.J. Michaels, D.M. Blackburn, E. Tabassi, and J.M. Bone, "Face Recognition Vendor Test 2002" FRVT.org
<http://www.frvt.org/DLs/FRVT_2002_Evaluation_Report.pdf>.

⁵ International Organization for Standardization, "Information technology -- JPEG 2000 image coding system: Conformance testing" ISO/IEC 15444-4:2004
<<http://www.iso.org/iso/en/CatalogueDetailPage.CatalogueDetail?CSNUMBER=39079&ICS1=35&ICS2=40&ICS3=&scopelist=>>.

⁶ International Organization for Standardization, "Standardization and related activities -- General vocabulary" ISO/IEC Guide 2:2004
<<http://www.iso.org/iso/en/CatalogueDetailPage.CatalogueDetail?CSNUMBER=39976>>.

⁷ National Institute of Standards and Technology, "Metrology for Information Technology (IT)" NISTIR 6025
<<http://www.itl.nist.gov/lab/nistirs/ir6025.htm>>.

About the National Science and Technology Council

The National Science and Technology Council (NSTC) was established by Executive Order on November 23, 1993. This Cabinet-level Council is the principal means within the executive branch to coordinate science and technology policy across the diverse entities that make up the Federal research and development enterprise. Chaired by the President, the membership of the NSTC is made up of the Vice President, the Director of the Office of Science and Technology Policy, Cabinet Secretaries and Agency Heads with significant science and technology responsibilities, and other White House officials.

A primary objective of the NSTC is the establishment of clear national goals for Federal science and technology investments in a broad array of areas spanning virtually all the mission areas of the executive branch. The Council prepares research and development strategies that are coordinated across Federal agencies to form investment packages aimed at accomplishing



Biometrics Testing and Statistics

multiple national goals. The work of the NSTC is organized under four primary committees; Science, Technology, Environment and Natural Resources and Homeland and National Security. Each of these committees oversees a number of sub-committees and interagency working groups focused on different aspects of science and technology and working to coordinate the various agencies across the federal government. Additional information is available at <http://ostp.gov/nstc>.

About the Subcommittee on Biometrics

Biometrics is a technology that is rapidly becoming a useful security, cost-savings and convenience tool for the Federal Government. Although the Federal Government is using the technology for many applications now, further development and assessment is required to improve the technology's utility. To address these issues, the Office of Science & Technology Policy (OSTP) created the NSTC Subcommittee on Biometrics, reporting to the National Science & Technology Council (NSTC) Committees on Technology and Homeland & National Security. Additional information is available at <http://www.biometricscatalog.org/NSTCSubcommittee>.

Subcommittee on Biometrics

Co-chair: Duane Blackburn (OSTP)
Co-chair: Chris Miles (DOJ)
Co-chair: Brad Wing (DHS)
Executive Secretary: Kim Shepard (FBI Contractor)

Department Leads

Mr. Jon Atkins (DOS)	Ms. Usha Karne (SSA)
Dr. Sankar Basu (NSF)	Dr. Michael King (IC)
Mr. Duane Blackburn (EOP)	Mr. Chris Miles (DOJ)
Ms. Zaida Candelario (Treasury)	Mr. David Temoshok (GSA)
Dr. Joseph Guzman (DoD)	Mr. Brad Wing (DHS)
Dr. Martin Herman (DOC)	Mr. Jim Zok (DOT)



Biometrics Testing and Statistics

Communications ICP Team

Champion: Duane Blackburn (OSTP)

Members & Support Staff:

Mr. Richard Bailey (NSA Contractor)

Mr. Jeffrey Dunn (NSA)

Ms. Valerie Lively (DHS S&T)

Mr. John Mayer-Splain (DHS US-VISIT Contractor)

Ms. Susan Sexton (FAA)

Ms. Kim Shepard (FBI Contractor)

Mr. Scott Swann (FBI)

Ms. Kimberly Weissman (DHS US-VISIT)

Mr. Brad Wing (DHS US-VISIT)

Mr. David Young (FAA)

Mr. Jim Zok (DOT)

Special Acknowledgements

The Communications ICP Team wishes to thank the following external contributors for their assistance in developing this document:

- FBI/OTD for authoring the 2003 "Biometrics 101" paper, which formed the basis of this document
- The Test and Evaluation ICP Team for reviewing the document and providing numerous helpful comments

Document Source

This document, and others developed by the NSTC Subcommittee on Biometrics, can be found at

<http://www.biometriccatalog.org/NSTCSubcommittee>.

