# IRI FieldShield
## PII / PHI Classification & Masking

# Product Overview



**Technical Summary, Samples, and Specifications**

# IRI
## Total Data Management

# Table of Contents

# Introduction

Consumers, patients, students and others expect their personally identifiable information (PII) and other sensitive data to remain protected from exposure, misappropriation and misuse. The companies and government agencies who collect, store, and/or process PII and other sensitive data are obliged to not only keep faith with these individuals, but assume the risks of financial liability or damage to their reputations for privacy law violations and data breaches.

The need to protect this data at risk -- at rest or in flight -- remains especially acute in traditional, transactional sources of data. Typically these are relational databases and structured files stored locally, or increasingly, in the cloud. They may also be used in testing and development.

One of the best ways to safeguard this data is by masking it through functions applied at the column level. This approach obfuscates data in secure and compliant ways that can nullify the effects of a data breach, differentiate access to original values, protect downstream applications from aforementioned risks, support auditing and still leave data usable for testing and analysis.

IRI FieldShield® is a robust data masking product for PII in databases and files. It was designed to help you adhere to both business rules and data privacy laws in a consistent, deterministic manner, either as a standalone solution or embedded in ETL, DevOps, or bespoke applications.

FieldShield is also a content-aware data loss prevention tool because it helps you find and apply the right 'shields' to only those fields that need masking. These field-level masks are also safer because even a breach of one field or encryption key still leaves remaining columns safe with their different protections. Contrast that to single password access to a full data source or silo.

# Operations

Most FieldShield users start with integrated data classification and discovery operations. They then create and run -- automatically or manually in an Eclipse IDE called IRI Workbench -- data masking tasks or batch jobs (for multiple sources at once), with or without using data classes or masking rules, and without relying on formally defined schema constraints.

Masking jobs are ultimately defined in FieldShield Control Language (.fcl) scripts, which use the same (SortCL) program syntax as IRI CoSort® for transformation, cleansing, and reporting, IRI NextForm® for data and database migration, and IRI RowGen® for synthesizing test data. FieldShield and all of the above SortCL spinoffs are included along with visual ETL and other big data management features in the IRI Voracity® platform.

FieldShield library functions documented in an SDK can also be used for dynamic data masking or unmasking, and to embed data protections into distributed applications or SQL calls.

For PII in semi-structured and unstructured document and file formats, including free form text, MF Office documents,.pdf's and images, and multiple NoSQL DBs, IRI DarkShield® can be used. Both FieldShield and DarkShield can also find and mask PII stored in Excel spreadsheets, but IRI CellShield® is purpose-built for operating on data directly within Excel. The data masking functions are the same between the tools, so the ciphertext should remain consistent across the enterprise if you use more than one IRI masking tool..

# Architecture

FieldShield works on – and can address DB and file sources in – LAN or cloud systems running Linux, Unix or Windows.

FieldShield uses IRI Workbench for client-side data discovery, and design of data-masking jobs, which serialize as portable text scripts that run from Workbench or any command line interface.

Many of the same scripts also run interchangeably in Hadoop MR2, Spark, Spark Stream, Storm or Tez for users of the big data edition of IRI Voracity.

FieldShield metadata is also fully supported in Workbench data models (EMF and XMI in Eclipse), and by Erwin Mapping Manager. Both offer graphical job creation, modification, and management.
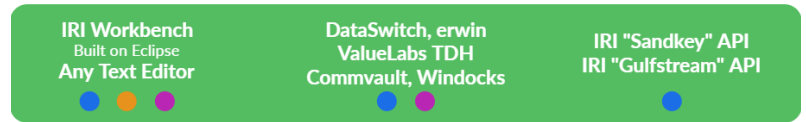
Separate Workbench wizards exist to discover PII, create metadata, and risk score the re-identifiability of data sets.

FieldShield masking functions, such as encryption, hashing, and redaction, can also be called into application programs written in C/C++, Java, and .NET through libraries in a separately available SDK.
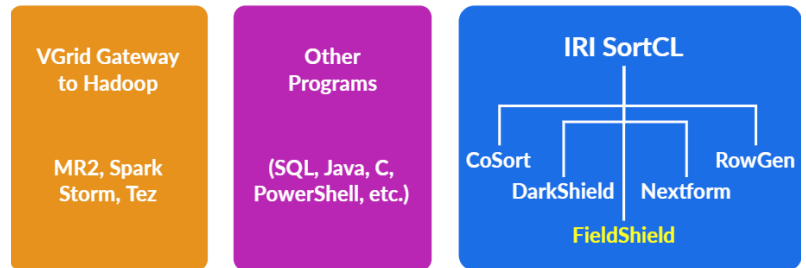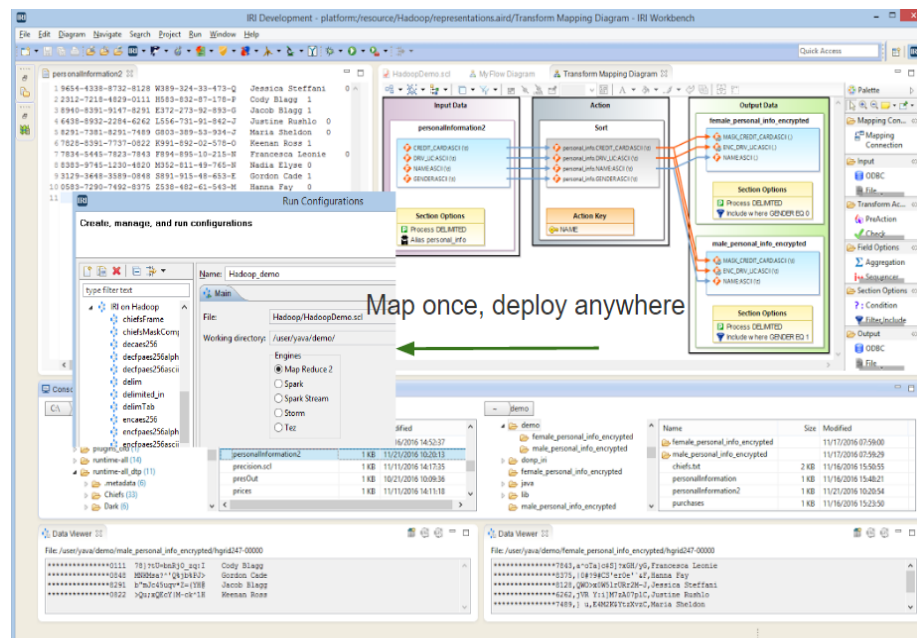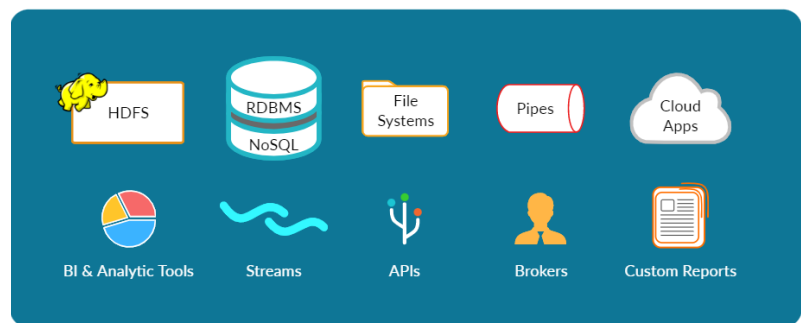
**IRI FieldShield**
PII / PHI Classification & Masking

### Job Designers

IRI Workbench
Built on Eclipse
**Any Text Editor**

**DataSwitch, erwin**
**ValueLabs TDH**
**Commvault, Windocks**

IRI "Sandkey" API
IRI "Gulfstream" API

### Runtime Engines

**VGrid Gateway to Hadoop**

MR2, Spark
Storm, Tez

**Other Programs**

(SQL, Java, C, PowerShell, etc.)

**IRI SortCL**

CoSort   RowGen
DarkShield   Nextform
**FieldShield**

### Sources & Targets

HDFS

RDBMS
NoSQL

File Systems

Pipes

Cloud Apps

BI & Analytic Tools

Streams

APIs

Brokers

Custom Reports

Map once, deploy anywhere

# PII Discovery

FieldShield's included data profiling and classification tools, pattern searches, forensic metadata discovery, and metadata management systems support data masking and governance by exposing the locations and attributes of PII before and after protection.

These tools are provided in IRI Workbench, the graphical IDE for FieldShield data discovery and masking jobs, built on Eclipse™. Included wizards support the discovery and definition of PII across disparate data sources, including DB tables, flat files, and dark (unstructured) data sources. Multiple reports from these search logs are produced, as are interactive dashboards.

These discovery processes can also support data cleansing, integration (ETL), subsetting, reformatting, test data generation, data wrangling, and reporting. Unique among data masking tools is this easy and guarded secret of performing such functions in conjunction with, or directly within, FieldShield masking and other metadata-compatible IRI data management operations.

# Data Classification & Masking

Define and manage enterprise-wide data class libraries, and use them to search and mask data across multiple sources at once. For more information, read this end-to-end how-to article::

https://www.iri.com/blog/data-protection/classify-mask-pii-in-databases-with-fieldshield/

or watch the video: www.youtube.com/watch?v=7nxp3TjsCqU&ab_channel=IRITheCoSortCo

# E-R Diagramming

Analyze the structure of, and relationships between, tables you select graphically in any connected database. For more information, see:

www.iri.com/blog/vldb-operations/er-diagrams-in-iri-workbench/

This wizard produces a diagram you can also customize, edit manually to change table attributes, and reproduce in different image formats for display and dissemination.

# Database & Flat-File Profiling

Through built-in profiling wizards, FieldShield users can produce statistics on the content and character of structured data sources, and search strings in them through matches to data in lookup files, Java RegEx patterns, and fuzzy-matching algorithms.

The data can also check the referential integrity of any (JDBC-connected) relational database.

For more information, see:
www.iri.com/blog/data-transformation2/database-profiling-in-iri-workbench/, and
www.iri.com/blog/iri/iri-workbench/flat-file-profiling/

# Unstructured Data Discovery

Because PII also exists in MS Office documents, .pdf, .rtf, and .xml files, logs and images, and other so called "dark data" repositories, it may be important to find, if not mask it. You can use the Dark Data Discovery wizard in IRI Workbench to define one or more patterns to search, and extract the values to DDF-defined flat files, along with forensic metadata about the sources discovered.



IRI DarkShield can be used for masking data in all of these formats, plus in database BLOB and CLOB columns where PII may be floating within unstructured document, text or image files. For database values exported to Excel, use either DarkShield, or the purpose-built spreadsheet-side masking tool, IRI CellShield EE.

Both DarkShield and CellShield EE can discover and de-identify data using the same data classes, and several of the same masking functions, used in FieldShield. They also share the same dark data discovery wizard for finding PII in multiple files and sheets simultaneously.

# Structured Data Discovery & Definition

IRI Workbench wizards connect to flat files and RDBs to find PII in multiple folders and schemas at once and facilitate the application of consistent masking rules. Searches can be performed within the data profiling wizards linked above, or from within fit-for-purpose wide-scan wizards.

Schema-wide pattern searches can identify columns matching supplied or custom regular expressions. Those results can be associated with data classes. In addition, schema-wide data class searches can find one or more data classes or groups which can activate multiple search methods at once. There is also a wizard for bulk masking across schema and DBs to save time and preserve referential integrity.



FieldShield users can also define and redefine column names, offsets, and data types automatically. They can also save and share the metadata in centralized data definition files and Git repositories. See: www.iri.com/blog/vldb-operations/table-filtering-iri-workbench
www.iri.com/blog/iri-iri-workbench/using-the-metadata-discovery-wizard
www.iri.com/blog/iri/iri-workbench/introduction-metadata-management-hub

# Static Data Masking Functions

The protection method you choose for each field depends on your need for security, reversibility, appearance, and speed. In addition to FieldShield's built-in functions, you can also write and call your own field functions at runtime.

FieldShield supports multiple masking functions, in multiple categories:

- multiple, NSA Suite B and FIPS-compliant encryption (and decryption) algorithms, including *format-preserving* encryption
- SHA-1 and SHA-2 hashing
- ASCII de-ID (bit scrambling)
- binary encoding and decoding
- blurring and bucketing (anonymization for HIPAA)
- random value generation or selection
- redaction (full or partial string masking)
- reversible and non-reversible pseudonymization
- filtering or omission (erasure for GDPR)
- literal or lookup value replacement
- byte shifting and (sub)string functions
- tokenization (for PCI DSS compliance)
- custom (user written and linked) routines

FieldShield uses state-of-the-art encryption algorithms that are very difficult to crack. And you can further strengthen encryption with random tokenization, hashing, and secure, or by changing encryption keys. Three encryption key management techniques are built in, and FieldShield is also compatible with the Microsoft Azure Key Vault and Alliance Key Manager from Townsend Security.

FieldShield data obfuscation techniques can also forever protect original values. A non-reversible method like redaction or filtering removes any computational basis for deriving the original value. You can also supply and link your own field-level protection functions for use at runtime with the others.

# Masking Job Wizards

Several automatic data masking job creation wizards are supplied in IRI Workbench. There are data class masking wizards for both tables in database schema and structured files in directories that leverage central data definitions and their associated search methods and masking rules. There are also older and simpler wizards that do not require data classification:

One works with a single data source, or multiple sources that share common fields, while the other wizard works with multiple tables with the same or different metadata, and allows you to apply cross-table protection rules to preserve referential integrity after protection occurs. Both wizards use the *New Field Rule Wizard*'s Selection page to help you create or use a saved protection rule.
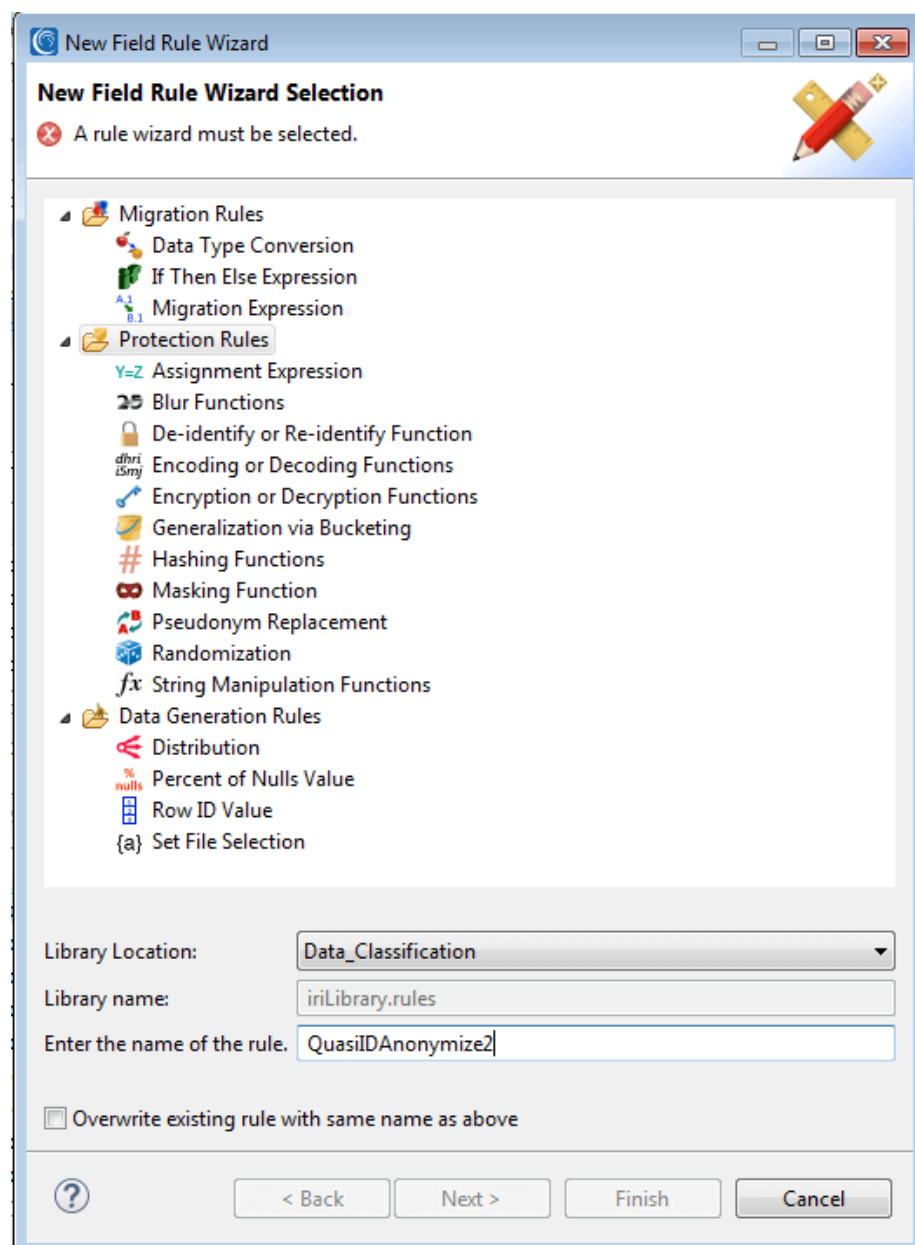
# Masking Rules

In the IRI Workbench [GUI for FieldShield](#), built on Eclipse, a variety of protection functions are available through the use of rule wizards that enable you to create simple or complex protection formatting for the script.

You can use the rule one time, or save it to a rule library for multiple or future use. The library is stored in the individual project. You can also create a rule library in each project, or create a rules project in which to store a single library for shared use.

When creating a new protection job with multiple tables and columns, use the Field Modification Rules dialog to create one of the many available types of protection.

When editing an existing job script, use the graphical Target Field Layout editor to create and add rules to the rule library, and then apply them to the job.
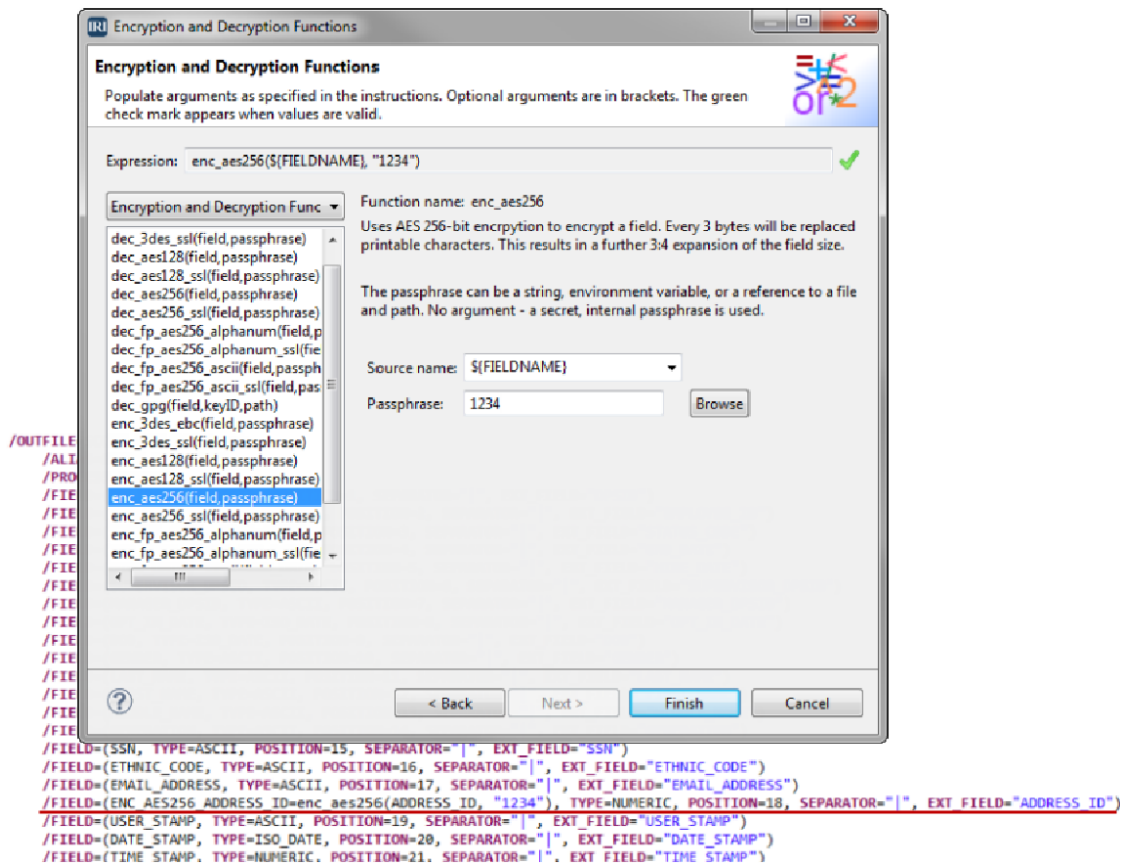
# Masking Function Dialogs

Most FieldShield functions can be applied visually in IRI Workbench dialogs, which appear in new job wizards, or when opened from the job script window or outline, or visual workflow.

# Encryption and Decryption

Encryption is one of the best ways to protect PII and primary account numbers (PANs) and other sensitive data because you use a specific key that makes the encrypted data generally unreadable. The data can only be decrypted using a matching decryption protocol and key.

FieldShield includes 3DES, AES, FIPS-compliant OpenSSL, and GPG encryption and decryption libraries. That dialog is shown below, but remember that you can also supply your own, custom encryption and decryption functions that you write and link.
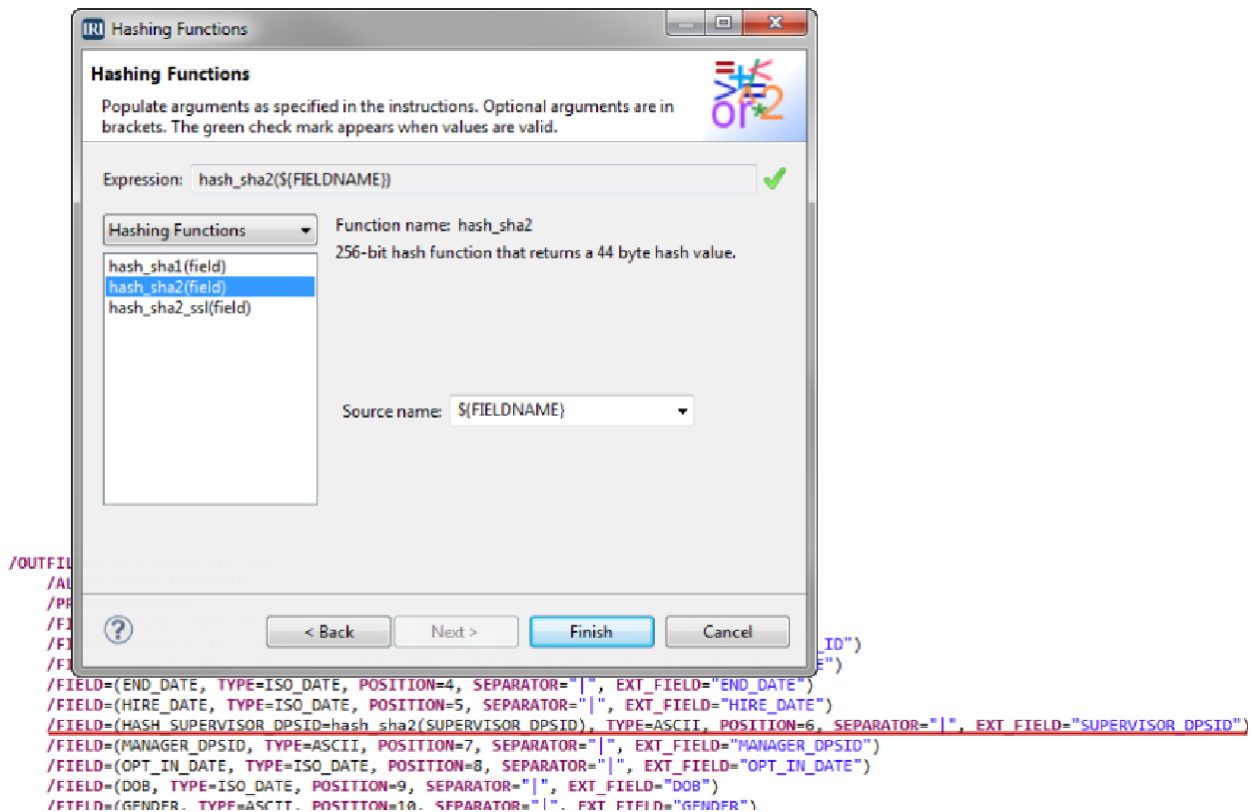


The most popular FieldShield function is 256-bit AES format-preserving encryption because the resulting ciphertext resembles the original field format and value.

FieldShield supports multiple encryption key management methods as well, to help authenticate users, and secure and rotate keys.

# Hashing

[Hashing](#) is a difficult-to-reverse data masking technique that converts a variable length "message" (e.g., someone's password) into an obfuscated, fixed-length, alphanumeric string. The purpose of the (Secure Hash Algorithm) SHA-2 field-level routine, for example, is to return a hash of the data string in a given field or column.
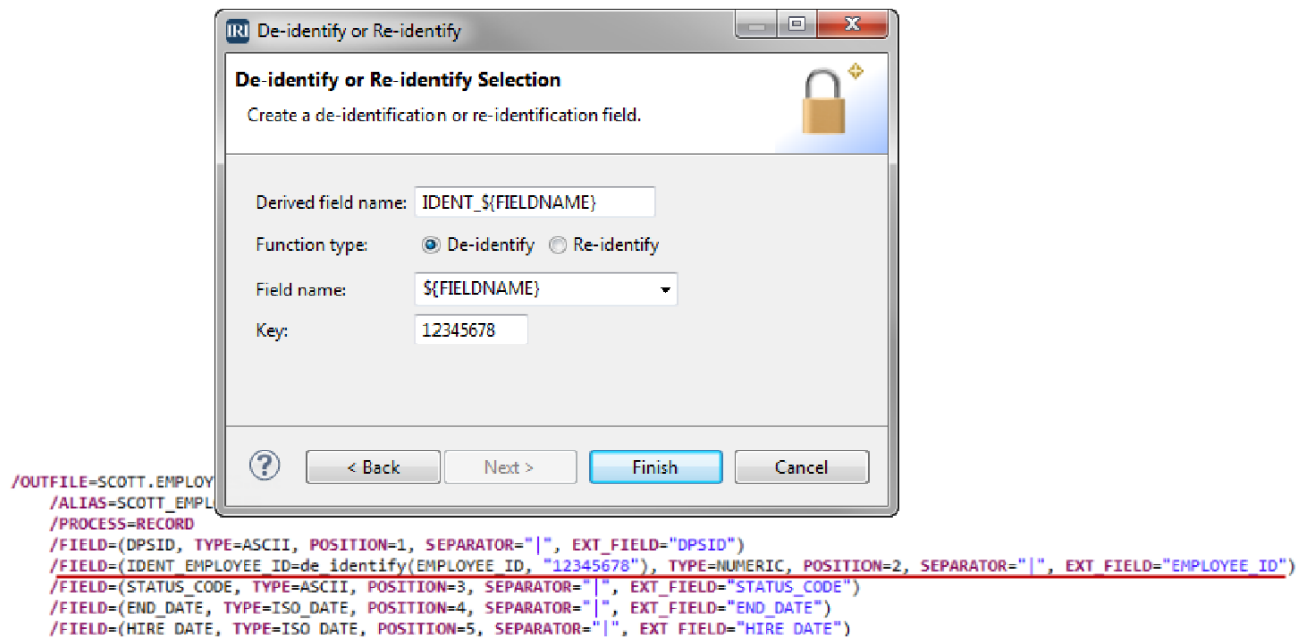


This is also known as a cryptographic hashing function. For any number of input bytes, the routine always returns a 32-byte hash code (in binary) that internally translates to a 44-byte, base64 value. Hashing is a useful technique for integrity checking. SHA-1, SHA-2, and SHA-2 OpenSSL are available.

# Bit Scrambling (ASCII De/Re-ID)

FieldShield can de- and re-identify PII, protected health information (PHI), and other sensitive field data in databases and files by [scrambling](#) data. This uses less computational overhead (and security) than encryption and decryption. At the same time, the de-identified field can look more realistic and still be recovered by using an encryption-style key.

The de-identification key used by FieldShield is a user-specified string, where the same string is used for subsequent re-identification. The hexadecimal range of the data contents to be de-identified must fall within the ASCII printable character range.
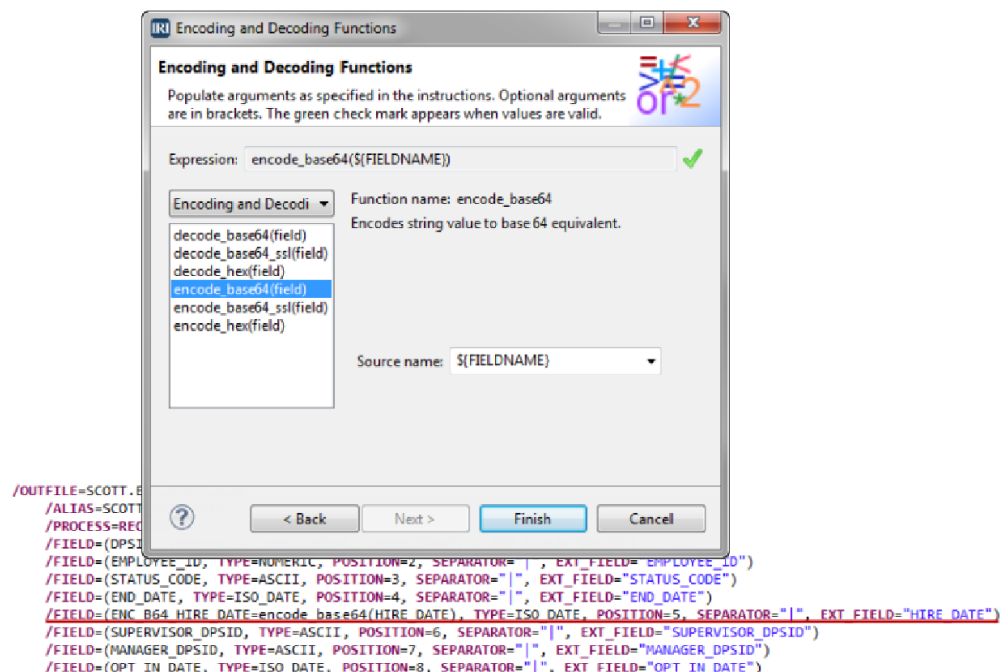
```
/OUTFILE=SCOTT.EMPLOY
   /ALIAS=SCOTT_EMPL
   /PROCESS=RECORD
   /FIELD=(DPSID, TYPE=ASCII, POSITION=1, SEPARATOR="|", EXT_FIELD="DPSID")
   /FIELD=(IDENT_EMPLOYEE_ID=de_identify(EMPLOYEE_ID, "12345678"), TYPE=NUMERIC, POSITION=2, SEPARATOR="|", EXT_FIELD="EMPLOYEE_ID")
   /FIELD=(STATUS_CODE, TYPE=ASCII, POSITION=3, SEPARATOR="|", EXT_FIELD="STATUS_CODE")
   /FIELD=(END_DATE, TYPE=ISO_DATE, POSITION=4, SEPARATOR="|", EXT_FIELD="END_DATE")
   /FIELD=(HIRE_DATE, TYPE=ISO_DATE, POSITION=5, SEPARATOR="|", EXT_FIELD="HIRE_DATE")
```

# Encoding and Decoding

Unlike encryption functions FieldShield supplies, encoding is not a particularly secure method of de-identification. While it does not require a passphrase or key to reverse, it is convenient and secure enough for some applications. Only the equivalent decoding function provided in FieldShield will reveal the original source values.

Use encoding routines when you need to encode byte data to be stored and transferred over media that are designed to deal with text data, such as email.

Three types of encoding convert the data in an ASCII string into its base64 equivalent value. An OpenSSL implementation of base 64 encoding is also available. The third type of encoding returns the hexadecimal representation of the byte input.



```
/OUTFILE=SCOTT.E
   /ALIAS=SCOTT
   /PROCESS=REC
   /FIELD=(DPSI
   /FIELD=(EMPLOYEE_ID, TYPE=NUMERIC, POSITION=2, SEPARATOR="|", EXT_FIELD="EMPLOYEE_ID")
   /FIELD=(STATUS_CODE, TYPE=ASCII, POSITION=3, SEPARATOR="|", EXT_FIELD="STATUS_CODE")
   /FIELD=(END_DATE, TYPE=ISO_DATE, POSITION=4, SEPARATOR="|", EXT_FIELD="END_DATE")
   /FIELD=(ENC_B64_HIRE_DATE=encode_base64(HIRE_DATE), TYPE=ISO_DATE, POSITION=5, SEPARATOR="|", EXT_FIELD="HIRE_DATE")
   /FIELD=(SUPERVISOR_DPSID, TYPE=ASCII, POSITION=6, SEPARATOR="|", EXT_FIELD="SUPERVISOR_DPSID")
   /FIELD=(MANAGER_DPSID, TYPE=ASCII, POSITION=7, SEPARATOR="|", EXT_FIELD="MANAGER_DPSID")
   /FIELD=(OPT_IN_DATE, TYPE=ISO_DATE, POSITION=8, SEPARATOR="|", EXT_FIELD="OPT_IN_DATE")
```
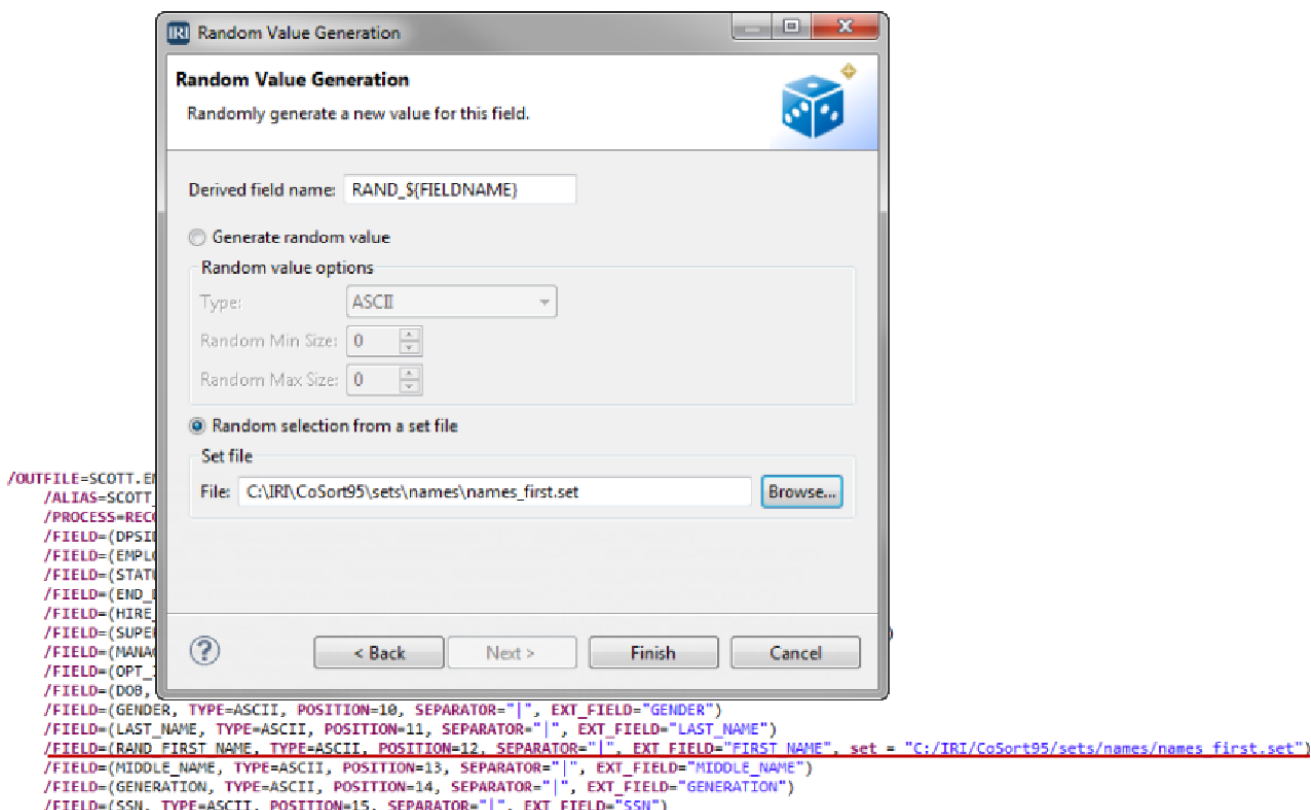
# Randomization

Replacing real values with random values is another secure, but non-reversible way to shield PII from disclosure, while maintaining the original structure of the input data. It can be utilized in two ways:

*Random data generation*, where random values are generated based on a data type with a random minimum and maximum size. Using random data, especially numbers, can make the protected field look somewhat real, while not requiring the use of any real data or protective overhead.

*Random selection from a set file, which* can use more realistic data from:

- an alpha-numeric list of values; the output values appear as they do in the set file.
- numeric values and ranges; the output values contain only the listed numbers, or a specified range of numbers.
- date values and ranges; only listed dates, or dates within a specific range, appear.



```
/OUTFILE=SCOTT.E
    /ALIAS=SCOTT
    /PROCESS=REC
    /FIELD=(DPSI
    /FIELD=(EMPL
    /FIELD=(STAT
    /FIELD=(END_
    /FIELD=(HIRE
    /FIELD=(SUPE
    /FIELD=(MANA
    /FIELD=(OPT_
    /FIELD=(DOB,
    /FIELD=(GENDER, TYPE=ASCII, POSITION=10, SEPARATOR="|", EXT_FIELD="GENDER")
    /FIELD=(LAST_NAME, TYPE=ASCII, POSITION=11, SEPARATOR="|", EXT_FIELD="LAST_NAME")
    /FIELD=(RAND_FIRST_NAME, TYPE=ASCII, POSITION=12, SEPARATOR="|", EXT_FIELD="FIRST_NAME", set = "C:/IRI/CoSort95/sets/names/names_first.set")
    /FIELD=(MIDDLE_NAME, TYPE=ASCII, POSITION=13, SEPARATOR="|", EXT_FIELD="MIDDLE_NAME")
    /FIELD=(GENERATION, TYPE=ASCII, POSITION=14, SEPARATOR="|", EXT_FIELD="GENERATION")
    /FIELD=(SSN, TYPE=ASCII, POSITION=15, SEPARATOR="|", EXT_FIELD="SSN")
```

# Redaction (String Masking)

Personally identifying information (PII) that is not needed for downstream processing or recovery should be redacted (masked) to prevent misuse. While DBAs can remove columns in tables, there are few alternatives for masking data in different ways across multiple databases and flat-file sources.

Only selected parts of the data should be exposed. At the same time, anonymization should preserve the original storage format and overall field appearance so that the platform structure or application need not be altered. Masking all but the last four digits of a credit card or social security number is a common requirement.

# Pseudonymization

Pseudonymization provides the level of realism you may require for your business. Two types of pseudonymization are available in FieldShield: recoverable and non-recoverable.

The recoverable function builds a lookup, or two-column set, file containing the original names from the input source, and substitution values randomly selected from that data to be paired with the original values. This effectively creates a shuffled list of real names so that there is no direct association between the original people and their attributes within other columns.

The non-recoverable function performs a random lookup from a supplied list, or single-column set file of names. You can choose among common western first, last or both names, as well as gender. This allows the names to appear real without compromising any one's identity.



IRI has also documented a method for maintaining unique pseudonyms and avoiding collisions as source values grow. See this article.

# Filtering or Removal

Removal is useful when sensitive data is not allowed in production or replicated data stores, or when it is just not needed in a future process step. You can easily remove columns or specific values from a database or file.



By removing data in this or other ways, FieldShield users can erase or delete source data in support of "Right to Be Forgotten" provisions of data privacy laws like the GDPR.

# Expression and String Functions

Numeric fields can be protected by applying mathematical calculation and trigonometric functions to the data. It is reversible only if you use the same expression for everything.

Bit-level functions can also manipulate column values. Available functions include string "find and replace", repositioning, trimming, and/or the use of INSTR information. Depending on the function and how you use it, the result of your manipulation may or may not be reversible or compliant.



```
/OUTFILE
    /ALI
    /PRO
    /FIE
    /FIE
    /FIE
    /FIE
    /FIE
    /FIE
    /FIE
    /FIE
    /FIE
    /FIE
    /FIE
    /FIE
    /FIE
    /FIE
/FIELD=(ETHNIC_CODE, TYPE=ASCII, POSITION=16, SEPARATOR="|", EXT_FIELD="ETHNIC_CODE")
/FIELD=(EMAIL_ADDRESS, TYPE=ASCII, POSITION=17, SEPARATOR="|", EXT_FIELD="EMAIL_ADDRESS")
/FIELD=(MASK_ADDRESS_ID=replace_chars(ADDRESS_ID, "+"), TYPE=NUMERIC, POSITION=18, SEPARATOR="|", EXT_FIELD="ADDRESS_ID")
/FIELD=(USER_STAMP, TYPE=ASCII, POSITION=19, SEPARATOR="|", EXT_FIELD="USER_STAMP")
/FIELD=(DATE_STAMP, TYPE=ISO_DATE, POSITION=20, SEPARATOR="|", EXT_FIELD="DATE_STAMP")
/FIELD=(TIME_STAMP, TYPE=NUMERIC, POSITION=21, SEPARATOR="|", EXT_FIELD="TIME_STAMP")
/FIELD=(CREATED_BY, TYPE=ASCII, POSITION=22, SEPARATOR="|", EXT_FIELD="CREATED_BY")
/FIELD=(CREATION_DATE, TYPE=ISO_DATE, POSITION=23, SEPARATOR="|", EXT_FIELD="CREATION_DATE")
/FIELD=(UPDATED_BY, TYPE=ASCII, POSITION=24, SEPARATOR="|", EXT_FIELD="UPDATED_BY")
/FIELD=(LAST_UPDATE_DATE, TYPE=ISO_DATE, POSITION=25, SEPARATOR="|", EXT_FIELD="LAST_UPDATE_DATE")
```
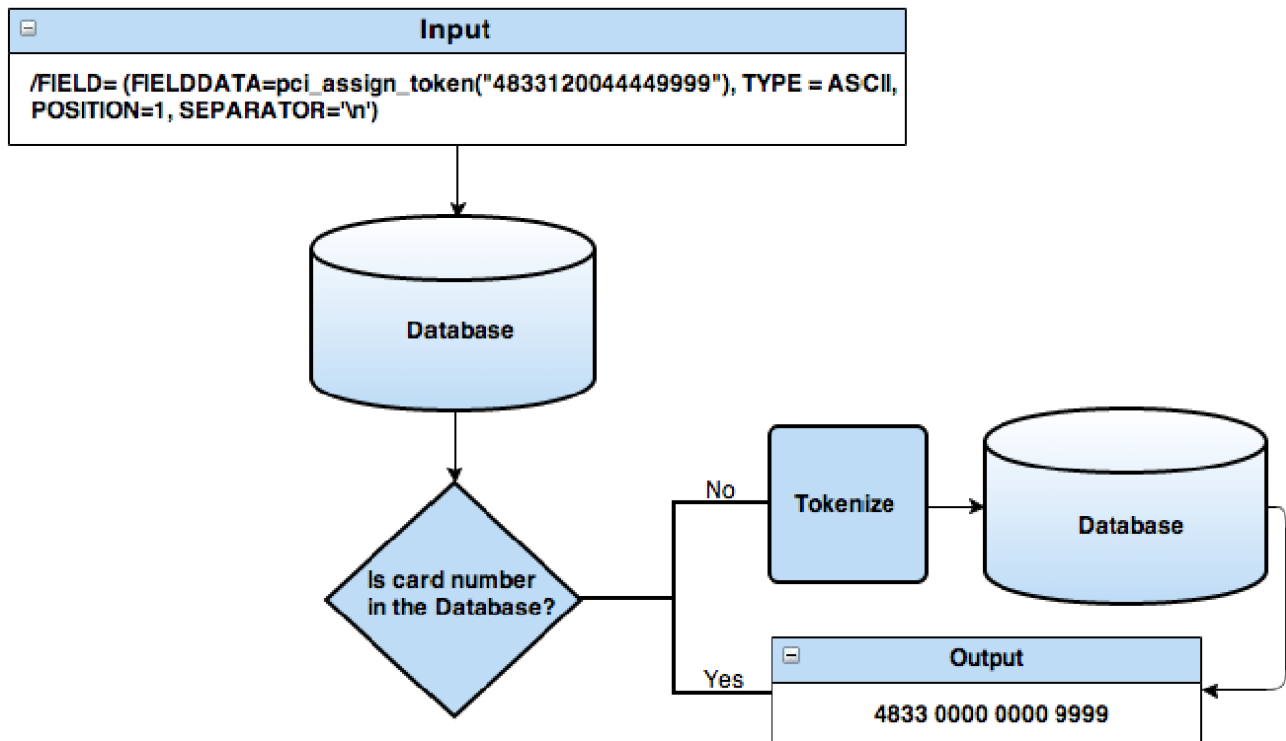
# Tokenization

[Tokenization](#) is another type of data masking operation, and is sometimes specified by the payment card industry (PCI) to be used in lieu of, or in conjunction with, Primary Account Number (PAN) encryption.
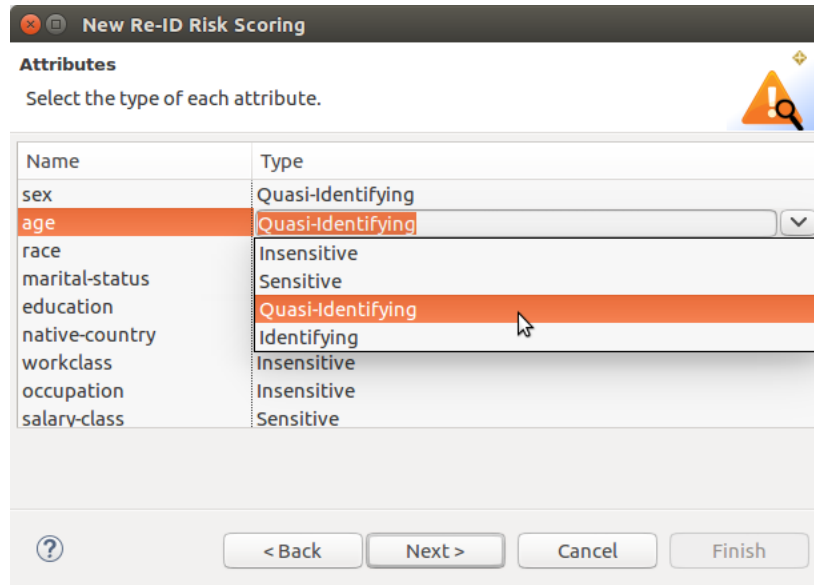
Unlike encryption where data is protected, but is still present, a token acts as a link to where the data is stored. Authorized users can tokenize credit card values that may also be treated with a Format-Preserving Encryption (FPE) function first, and possibly de-tokenized later. The diagram below depicts the tokenization process.



Contact [fieldshield@iri.com](mailto:fieldshield@iri.com) if you require another implementation of tokenization, or if you need help with or documentation on the use of FieldShield control language in searching selections, protecting, or reformatting data.
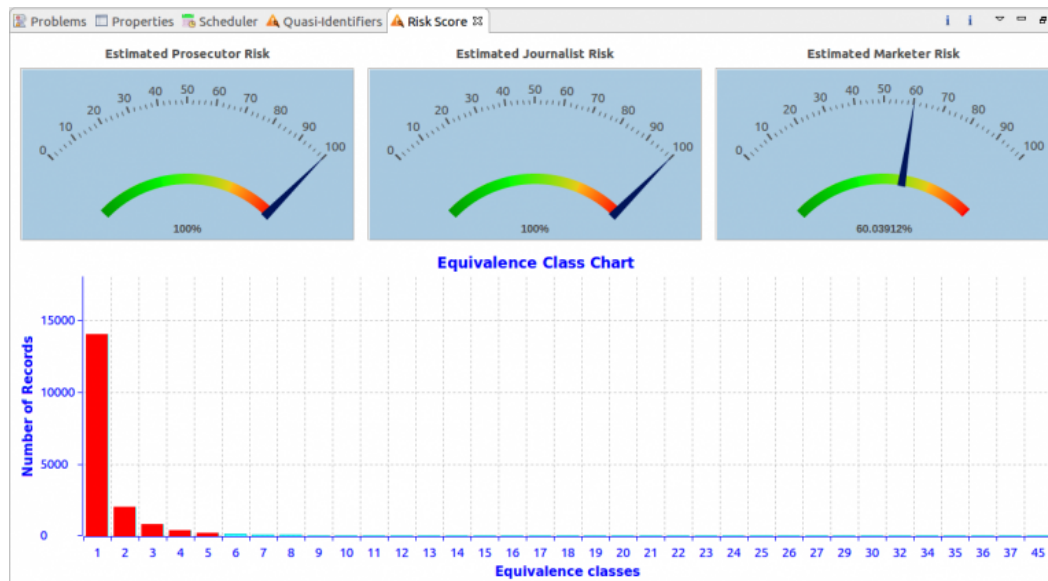
# Re-ID Risk Scoring

To comply with the Health Insurance Portability and Accountability Act (HIPAA), you must de-identify both *key* identifiers and *quasi*-identifiers. Key identifiers are unique PHI values like name and Social Security Number, while quasi-identifiers are less unique attributes like age, race, state, gender, and occupation which can be used in tandem to identify unique records.
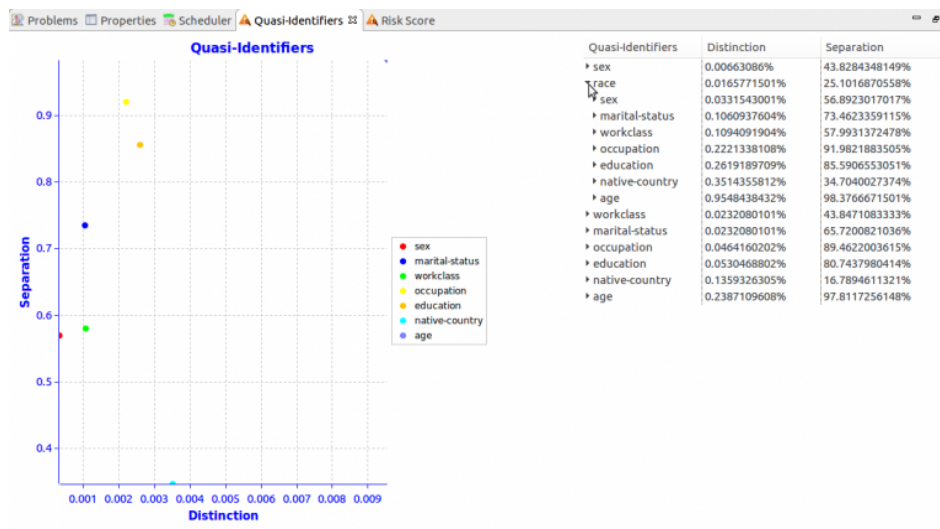


The Re-ID Risk Scoring wizard in IRI Workbench produces detailed and visual reports that statistically measure the risk of re-identification in three modes of attack and the number of records in each equivalence class:



Another chart provides an interactive look at the different combinations of quasi-identifiers and their separation and distinction values to further assess their capacity to re-identify a record:

| Quasi-Identifiers | Distinction | Separation |
|---|---|---|
| ▸ sex | 0.00663086% | 43.8284348149% |
| race | 0.0165771501% | 25.1016870558% |
| ▸ sex | 0.0331543001% | 56.8923017017% |
| ▸ marital-status | 0.1060937604% | 73.4623359115% |
| ▸ workclass | 0.1094091904% | 57.9931372478% |
| ▸ occupation | 0.2221338108% | 91.9821883505% |
| ▸ education | 0.2619189709% | 85.5906553051% |
| ▸ native-country | 0.3514355812% | 34.7040027374% |
| ▸ age | 0.9548438432% | 98.3766671501% |
| ▸ workclass | 0.0232080101% | 43.8471083333% |
| ▸ marital-status | 0.0232080101% | 65.7200821036% |
| ▸ occupation | 0.0464160202% | 89.4622003615% |
| ▸ education | 0.0530468802% | 80.7437980414% |
| ▸ native-country | 0.1359326305% | 16.7894611321% |
| ▸ age | 0.2387109608% | 97.8117256148% |

In addition to those interactive graphs which can be saved in different image formats, the FieldShield re-ID risk determination report provides detailed descriptions of the metrics:



## Uniqueness Estimates

The uniqueness of individuals in a population is often used as a measure of re-identification risk. Another convenient way to assess the uniqueness of records in a dataset is to assess their uniqueness relative to a statistical sample of the superset to which the dataset conforms. Four different statistical methods have been identified as the best measures of relative uniqueness, and their scores follow. The percentage represents the conditional probability that a record in the superset is unique if it is unique in the original data set.

> Number of uniques in sample: 0
> Percentage of uniques relative to the USA superset:
> Dankar: 2.97%
> Pitman: 2.97%
> slide negative binomial (SNB): 0.00%
> Zayatz: 21.62%

Note that unique records within the dataset are at the highest risk of being re-identified. However, a unique record within the dataset may not automatically be at risk unless it can joined to a unique record within the superset that contains the record's identity. Conversely, a record with a higher equivalence class (lower re-identification risk

After reviewing the risk scoring report in consultation with a qualified statistician whom IRI can also refer, you can create additional FieldShield jobs that generalize or blur one or more of the quasi-identifiers so they remain useful for research or marketing purposes, but less likely to lead to re-identification. See the section on generalization. After that, you can easily re-score re-generalized data sets from the attribute model you created in the first pass through the wizard.

For more information, see:
www.iri.com/blog/data-protection/hipaa-re-id-risk-scoring/

# Anonymization (Blurring and Bucketing)

Quasi-identifying numeric values like age and date of birth, as well as descriptors like occupation and marital status, can all be used to re-identify people if there are enough of these attributes in the data set and/or they can be joined to a superset population with similar values.

For this reason, FieldShield jobs can apply one or more additional techniques to further anonymize or obfuscate the data, while still keeping it accurate enough for research or marketing purposes.

A numeric "smart blurring" function allows specific random values or values within defined ranges to be applied to ages and dates. This is also referred to as random noise.

In the example below however, specific ages are bucketed into decade groups, multiple marital status attributes are combined into two broader categories in a defined condition, educational attainments are simplified through a new set lookup file, and all occupations were explicitly redacted in place.



The new result set can now be re-run through the risk scoring wizard to produce another determination of re-identification risk based on now less distinct quasi-identifying attributes.

# Software Development Kit / DDM

The software development kit ([SDK](#)) for dynamic data masking (DDM) or unmasking within applications uses APIs written in Java and .NET.

The FieldShield library is a group of classes that provides several ways of protecting or altering data by using encryption/decryption processes, hashing, masking, and base64 to hex conversion. These functions are compatible with those used in the standalone, static data masking executable. It is therefore possible to encrypt data at rest in a database in a standalone FieldShield operation, and selectively mask it in transit from an application calling a corresponding decryption function in the SDK.

FieldShield SDK encryption and decryption algorithms include:

- AES-256 (Advanced Encryption Standard); the state-of-the-art and most powerful of the encryption algorithms.
- FP-ASCII (Format Preserved for American Standard Code for Information Interchange); the original format data is preserved but still protected with AES-256 bit encryption.
- FP- ALPHANUM; another version of the Format Preserving algorithm above, in which it can be set to number mode only or alphanumeric for data preservation.

# FieldShield Library Process

Each function in the FieldShield library flows in a defined manner. The figure below shows the sequence of methods for each encryption and decryption usage. Once the password is set for a given object, it does not need to be set each time the encryption object is used.



The other classes follow a similar flow in usage, varying slightly. The steps are defined thus:

1. **Creation** - Create an instance of the desired class. Most languages use the keyword new.

2. **Begin** - After creating an instance from the class, use this instance to call the begin() method. This function will initiate an encryption or decryption process.

3. **SetPass** - Next, the pass phase must be set by calling the setPass method. This method requires a String to be passed in as a parameter.

4. **DoTransform** - Once the passphrase has been set, the core method, doTransform(Data), is called to encrypt or decrypt. A String is passed as a parameter, which is then encrypted or decrypted using the algorithm chosen. For the alphanumeric class, a numeric doTransform method exists to be used on numerals. The method will return the data that has been encrypted or decrypted with the appropriate type.

5. **Destruct** - Finally, call the method destruct() to release the object from the occupied space in memory. Since it uses the DLL file, it is necessary to destruct the object.

6.

# Code Samples

Java Example

```
enc_aes256 obj = new enc_aes256();
iResult = obj.begin(obj.getptr());
obj.setPass(obj.ptr, "password");
String enc = obj.doTransform(obj.getptr(), "String to perform encryption on");
obj.destruct(obj.getptr());
```

.NET Example in C#

```
enc_aes256 encryptAES256 = new enc_aes256();
iResult = encryptAES256.begin();
encryptAES256.setPass("password");
string test1 = encryptAES256.doTransform("String to perform encryption on");
encryptAES256.end();
```

# Technical Specifications

FieldShield is a program for protecting Personally-Identifying Information (PII), Protected Health Information (PHI), and other sensitive data using highly efficient, targeted data masking functions according to business rules. It can nullify the effect of data breaches by protecting PII at the field level across multiple data sources and facilitate compliance with U.S. and international data privacy laws.

FieldShield can also provide file-level protection with one encryption algorithm and a key. However, most FieldShield users apply multiple functions to specific elements; i.e., columns in database tables and fields in files that are intended for specific uses and recipients.

To safeguard sensitive data, FieldShield offers four significant features that enable you to do the following in any combination:

- Choose one or more protection (e.g., encryption) functions in combination
- Apply theses functions at the field (column) or record (row) level
- Make secure data look real
- Generate secure output targets in one pass through the input data

## Installation

- Distributes via Internet or user-specified media
- Loads in under two minutes
- Uses menu-driven setup and configuration utility

## Invocation

- Command line (including pipe sequences), shell commands and batch scripts
- IRI Workbench – Eclipse GUI for Windows (JRE required)
- Application calling programs as a standalone executable, subroutine or coroutine call, with or without additional exit routines

## Ease of Use

- Processes data using record layouts and SQL-like field descriptions within applications or centralized data definition files (DDF)
- Provides on-line help, pre-runtime application validation, and runtime errors
- Leverages centralized application and file layout definitions (metadata repositories)
- Reports problems to standard error when invoked from a program, or to an error log
- Runs silently or with verbose messaging without user intervention
- Allows user control over the amount of informational output produced
- Generates a query-ready XML audit log for data forensics and privacy compliance
- Describes commands and options through on-line documentation
- Fully supports job design, execution, and modification in familiar Eclipse GUI
- Easy-to-use interfaces preclude the need for training classes (advanced training is available in Florida or at user sites)
- Phone, web, and email support available directly from the product developers
- Local language support is available from more than 30 international offices

## Functions

- Data discovery and classification through multiple search and profiling wizards
- Multiple encryption via AES-256 routines (e.g., FPE_ASCII / FPE_ALPHANUM) or a user-supplied routine
- GPG/PGP encryption and decryption with GPG key ring management support
- ASCII de-identification, encoding, hashing, and randomization
- Reversible and non-reversible pseudonymization
- Redaction (partial or full string masking) or conditional value removal
- Anonymization via generalization, blurring, or custom obfuscation
- Re-ID risk determination and reporting through a fit-for-purpose scoring wizard
- Audit logging to verify compliance, and statistical reporting for performance review

# Data Sources

## Sources

Local pipes via stdin, plus any number of structured (COBOL, CSV, LDIF, fixed, etc.), and *some* semi-structured (flat JSON and XML, Excel, ASN.1 CDR) files on-premise *or* in Amazon S3, Azure Blob or GCP storage buckets. Also supported are named pipes (on UNIX or Linux), and ODBC-connected relational databases on premise or in the cloud. HDFS, MongoDB, HIVE, MQTT, Kafka, HTTP/S and FTP/S are also accessible through URL specifications.

## Input Size

No limit

## Record/Row Length

1 ≤ fixed-length ≤ 65,535 bytes
0 ≤ variable-length ≤ 65,535 bytes

## Program

Optional, user-written routines for reading special sources

## Data Types

The following categories of data types are supported:

| Character strings | ASCII, EBCDIC, ASCII in EBCDIC order, days of the week, months of the year, etc. |
|---|---|
| Numerics | ASCII-Numeric, MF COBOL types, RM COBOL types, integers, floats, doubles, zoned decimal, etc |
| Timestamps | American, European, ISO, Japanese, and current time, date, and timestamp. |
| Multi-byte | Double-byte types such as Shift JIS, EHANGUL, and Big 5. |

# Targets

## Terminal

stdout

## Table

ODBC- and JDBC-connected databases on premise or in the cloud

## File

User-selected name(s) or device(s), including input file overwrite and create or append logic, plus named and unnamed pipes. Target files can persist in the LAN or cloud storage buckets (specifically, in Amazon S3, Azure Blob, and Google Cloud Platform)..

## Both

Standard output and file to observe and abort unintended executions

## URLs

HDFS, MongoDB, HIVE, MQTT, Kafka, and FTP/S

## Report

Customized detail report with header, footer, page count, labeling, and other formatting features

## Program

Optional, user-written routines for writing to special targets

## Operating Optional Modes

- standalone or embedded batch execution
- user program integration

# Platforms

- AIX 6.1 and above on IBM Power architecture
- HP-UX 11.31b on ia64 architecture
- Linux kernel 2.6 and above on Intel x86 & x64 architecture
- Linux kernel 3.6 and above on ARMv6 and above architecture
- macOS
- Solaris 5.8 and above on Sun SPARC architecture
- Solaris 10 on Intel x64 architecture
- Windows XP and above on Intel x86 & x64 architecture
- Windows Server 2003 and above on Intel x86 & x64 architecture

# Licensing Information

IRI and its representatives around the world license FieldShield for leased or perpetual use on individual computer systems. Maintenance (technical support and site-specific software updates) services are provided free of charge during the first year after installation. Subsequent annual maintenance is usually offered at 20% of the base license fee, and 24/7 technical support is available.

Discounts are available for multiple copies of FieldShield at the same site, distributed anonymization, and for runtime integration and redistribution (ISVs only). IRI is generous with credit for hardware upgrades and migrations, and provides for no- or low-cost failover (disaster recovery) licenses.

U.S. educational and 501c(3) non-profit institutions qualify for a 10% license fee discount, and government agencies can buy FieldShield from prime contractors on GSA schedule.

FieldShield is also a product member of the [IRI Data Protector Suite](#) and an included component in the [IRI Voracity](#) total data management platform (subscription). Therefore, its functionality can be licensed in a discounted bundle with other IRI data-centric security software.


# Professional Services

An IRI professional services engagement allows you to leverage more than 100 collective years of IT and integrated data-handling experience. Examples of IRI professional services engagements involve:

- *Big data preparation* – packaging, protecting, and provisioning structured and unstructured data sets for analytic/BI, database (DB), and ETL operations
- *Data masking* – profiling, de-identification, encryption, tokenization and other services to aid data loss prevention, data governance, and data privacy law compliance efforts
- *Risk scoring and certification* -- partner statistical analysis and legal advice pertaining to data re-identification, HIPAA compliance, insurance coverage, and breach defense
- *Data replication and federation* – acquiring, re-mapping, and creating views
- *DB migration* – mapping table data and relations to new versions/platforms
- *Data conversion* – reformatting LDIF, XML, and COBOL index files (e.g., Vision, MF-ISAM), multi-byte character sets, most mainframe data types, and endian states
- *Master data management* – value discovery, format definition, validation and enrichment, security, and lineage
- *Program replacement* – translating cryptic and inefficient SQL, 3GL, ETL, legacy sort, and shell procedures into IRI's simple, portable, 4GL text scripts
- *Test data management* – end-to-end services from DevOps needs definition through data generation and target persistence (resembling production data)

**INNOVATIVE ROUTINES INTERNATIONAL (IRI), INC.**

Suite 303, Atlantis Center
2194 Highway A1A
Melbourne, FL 323937-4932 USA
Phone | +1.321.777.8889
www.iri.com/fieldshield