

NIST SPECIAL PUBLICATION 1800-39

Data Classification Practices

William Newhouse

Murugiah Souppaya*

National Institute of Standards and Technology
Gaithersburg, Maryland

John Kent

Kenneth Sandlin*

Ryan Williams*

The MITRE Corporation
McLean, Virginia

Karen Kent

Trusted Cyber Annex

Mark Evans

Jimmy Katz*

ActiveNav

John Dombroski

Harmeet Singh

IBM

Neville Jones

Helen Farrell*

Janusnet

Pablo Blasco

Jane Gilbert

D'Nan Godfrey

Matt Jochim

Rich Johnson

Ludmila Rinaudo

Gina Scinta

Thales Trusted Cyber Technologies

Patrick Greer

Wilson Patton

Jason White*

Trellix

*Former employee; all work for this publication was done while at listed organization.

February 2026

INITIAL PUBLIC DRAFT

This publication is available free of charge from

<https://www.nccoe.nist.gov/data-classification>



1 **DISCLAIMER**

2 Certain commercial entities, equipment, products, or materials may be identified by name or company
3 logo or other insignia in order to acknowledge their participation in this collaboration or to describe an
4 experimental procedure or concept adequately. Such identification is not intended to imply special sta-
5 tus or relationship with NIST or recommendation or endorsement by NIST or NCCoE; neither is it in-
6 tended to imply that the entities, equipment, products, or materials are necessarily the best available
7 for the purpose.

8 While NIST and the NCCoE address goals of improving management of cybersecurity and privacy risk
9 through outreach and application of standards and best practices, it is the stakeholder’s responsibility to
10 fully perform a risk assessment to include the current threat, vulnerabilities, likelihood of a compromise,
11 and the impact should the threat be realized before adopting cybersecurity measures such as this
12 recommendation.

13 National Institute of Standards and Technology Special Publication 1800-39, Natl. Inst. Stand. Technol.
14 Spec. Publ. 1800-39 60 pages, (February 2026), CODEN: NSPUE2

15 **Author ORCID iDs**

16 William Newhouse: 0000-0002-4873-7648

17 Murugiah Souppaya: 0000-0002-8055-8527

18 John Kent: 0009-0001-9989-2277

19 Ken Sandlin: 0009-0000-3757-4858

20 Karen Kent: 0000-0001-6334-9486

21 **FEEDBACK**

22 You can view or download the guide at the [NCCoE Data Classification Practices project page](#).

23 With this initial public draft, NIST NCCoE is now asking for feedback on this initial public draft.

24 Comments on this publication may be submitted to: data-nccoe@nist.gov.

25 Public comment period: February 12, 2026 - March 30, 2026

26 All comments are subject to release under the Freedom of Information Act.

27 NIST NCCoE is particularly interested in your feedback on the following questions:

- 28 1. Does this guide influence your plans to use of data classification practices?
29 2. How well do the data classification practices in this guide relate to your organization’s practices
30 or needs?
31 3. Are there significant gaps this guide should address?
32 4. What would you like to see added to this guide?
33 5. Do you have suggestions for future NCCoE data centric cybersecurity or privacy projects?
34 6. What could the NCCoE build as an example demonstration using commercially available
35 technologies that would help you take action to mitigate some aspects of your cybersecurity and
36 privacy risks?

37

38

National Cybersecurity Center of Excellence

39

National Institute of Standards and Technology

40

100 Bureau Drive

41

Mailstop 2002

42

Gaithersburg, MD 20899

43

Email: nccoe@nist.gov

44 **NATIONAL CYBERSECURITY CENTER OF EXCELLENCE**

45 The National Cybersecurity Center of Excellence (NCCoE), a part of the National Institute of Standards
46 and Technology (NIST), is a collaborative hub where industry organizations, government agencies, and
47 academic institutions work together to address businesses' most pressing cybersecurity issues. This
48 public-private partnership enables the creation of practical cybersecurity solutions for specific
49 industries, as well as for broad, cross-sector technology challenges. Through consortia under
50 Cooperative Research and Development Agreements (CRADAs), including technology partners—from
51 Fortune 50 market leaders to smaller companies specializing in information technology security—the
52 NCCoE applies standards and best practices to develop modular, adaptable example cybersecurity
53 solutions using commercially available technology. The NCCoE documents these example solutions in
54 the NIST Special Publication 1800 series, which maps capabilities to the NIST Cybersecurity Framework
55 and details the steps needed for another entity to re-create the example solution. The NCCoE was
56 established in 2012 by NIST in partnership with the State of Maryland and Montgomery County,
57 Maryland.

58 To learn more about the NCCoE, visit <https://www.nccoe.nist.gov/>. To learn more about NIST, visit
59 <https://www.nist.gov>.

60 **NIST CYBERSECURITY PRACTICE GUIDES**

61 NIST Cybersecurity Practice Guides (Special Publication 1800 series) target specific cybersecurity
62 challenges in the public and private sectors. They are practical, user-friendly guides that facilitate the
63 adoption of standards-based approaches to cybersecurity. They show members of the information
64 security community how to implement example solutions that help them align with relevant standards
65 and best practices, and provide users with the materials lists, configuration files, and other information
66 they need to implement a similar approach.

67 The documents in this series describe example implementations of cybersecurity practices that
68 businesses and other organizations may voluntarily adopt. These documents do not describe regulations
69 or mandatory practices, nor do they carry statutory authority.

70 **ABSTRACT**

71 This guide demonstrates how organizations can discover, identify and label unstructured data using data
72 classification practices. Performing Data Classification Practices allows an organization to know its data
73 and apply technologies that minimize the risk of valuable or sensitive data being lost or mismanaged.
74 Data Classification Practices prepare an organization for the use of emerging security measures—
75 including Zero Trust Architecture, quantum-safe cryptography, and AI model training that requires
76 labeled data. This 1800-series NIST publication documents how the NCCoE and its collaborators created
77 a synthetic dataset and used commercially available data classification tools to discover, identify and
78 label unstructured data.

79 **KEYWORDS**

80 Data classification practices; synthetic unstructured data; data pillar of zero trust; data types, data
81 labels.

82 **ACKNOWLEDGMENTS**

83 We are grateful to the following individuals from our project’s technology collaborators for their
84 generous contributions of expertise and time.

- 85 ▪ Adobe: Steve Gottwals
- 86 ▪ GitLab: Joel Krooswyk
- 87 ▪ JPMorgan Chase & Co.: Timothy Brophy*, Lauren Brown, Benjamin Flatgard
- 88 ▪ MITRE: Jason Ajmo, Lauren Swan, Spike Dog*
- 89 ▪ NIST: Cherilyn Pascoe, Julie Chua
- 90 ▪ Seqrite/Quick Heal: Dhruvi Desai, Nachiket Karguppikar, Vinaya Sathyanarayana*
- 91 ▪ Virtru: Will Ackerly, Cassandra Zimmerman

92 * Former employee; all work for this publication was done while at listed organization.

93 We are also grateful to the experts from JPMorgan Chase, Microsoft, Morgan Stanley, NATO, NIST, and
94 Varonis who presented at an October 2019 NCCoE hosted Data Classification workshop. Their
95 contributions led to the establishment of the NCCoE Data Classification project.

96 We wish to thank Bill Brunt for his interest in our project and his willingness to share his data
97 classification practices expertise at the start of our project.

98 The authors are grateful to all who reviewed and provided feedback on this publication.

99 The Technology Collaborators who participated in this build submitted their capabilities in response to a
100 notice in the Federal Register. Respondents with relevant capabilities or product components were
101 invited to sign a Cooperative Research and Development Agreement (CRADA) with NIST, enabling their
102 participation in a consortium that built this publication’s example demonstrations.

103 **DOCUMENT CONVENTIONS**

104 The terms “shall” and “shall not” indicate requirements to be followed strictly to conform to the
105 publication and from which no deviation is permitted. The terms “should” and “should not” indicate that
106 among several possibilities, one is recommended as particularly suitable without mentioning or
107 excluding others, or that a certain course of action is preferred but not necessarily required, or that (in
108 the negative form) a certain possibility or course of action is discouraged but not prohibited. The terms
109 “may” and “need not” indicate a course of action permissible within the limits of the publication. The
110 terms “can” and “cannot” indicate a possibility and capability, whether material, physical, or causal.

111 **CALL FOR PATENT CLAIMS**

112 This public review includes a call for information on essential patent claims (claims whose use would be
113 required for compliance with the guidance or requirements in this Information Technology Laboratory
114 (ITL) draft publication). Such guidance and/or requirements may be directly stated in this ITL Publication
115 or by reference to another publication. This call also includes disclosure, where known, of the existence
116 of pending U.S. or foreign patent applications relating to this ITL draft publication and of any relevant
117 unexpired U.S. or foreign patents.

118 ITL may require from the patent holder, or a party authorized to make assurances on its behalf, in
119 written or electronic form, either:

120 a) assurance in the form of a general disclaimer to the effect that such party does not hold and does not
121 currently intend holding any essential patent claim(s); or

122 b) assurance that a license to such essential patent claim(s) will be made available to applicants desiring
123 to utilize the license for the purpose of complying with the guidance or requirements in this ITL draft
124 publication either:

- 125 1. under reasonable terms and conditions that are demonstrably free of any unfair
126 discrimination; or
- 127 2. without compensation and under reasonable terms and conditions that are
128 demonstrably free of any unfair discrimination.

129 Such assurance shall indicate that the patent holder (or third party authorized to make assurances on its
130 behalf) will include in any documents transferring ownership of patents subject to the assurance,
131 provisions sufficient to ensure that the commitments in the assurance are binding on the transferee,
132 and that the transferee will similarly include appropriate provisions in the event of future transfers with
133 the goal of binding each successor-in-interest.

134 The assurance shall also indicate that it is intended to be binding on successors-in-interest regardless of
135 whether such provisions are included in the relevant transfer documents.

136 Such statements should be addressed to: data-nccoe@nist.gov.

137 **Contents**

138 **1 Overview 1**

139 1.1 Challenge1

140 1.2 Audience1

141 1.3 Scope2

142 1.4 Structure of This Guide2

143 **2 Project Overview 3**

144 2.1 Motivation for the Project3

145 2.2 Challenges in Implementing Data Classification Practices for Unstructured Data3

146 2.3 Project Approach4

147 **3 Synthetic Data..... 5**

148 3.1 Synthetic Data Characteristics5

149 3.2 Synthetic Data Source5

150 3.3 Unstructured Synthetic Data Files5

151 **4 Demonstrations 9**

152 4.1 Unstructured Data Classification Practice Demonstrations9

153 4.2 Lab Demonstration Environment9

154 4.3 Unstructured Data Classification Practice Demonstration Workflow10

155 4.4 Electronic Mail Message Demonstration17

156 4.5 Tool Summary21

157 **5 Findings and Insights..... 25**

158 **Appendix A List of Acronyms 26**

159 **Appendix B Glossary 27**

160 **Appendix C References 28**

161 **Appendix D Synthetic Data Creation Steps 29**

162 **Appendix E Lab Implementation Details..... 38**

163 **List of Figures**

164 **Figure 3-1 Specific Example of Synthetic Unstructured Data for Each Data Element 8**

165 **Figure 4-1 High Level Overview of Lab Environment for Each Data Classification Tool 10**

166 **Figure 4-2 Unstructured Data Classification Practice Demonstration Workflow 12**

167 **Figure 4-3 Email Message Demonstration Workflow 19**

168 **Figure D-1 Synthetic Data Template Example 30**

169 **Figure D-2 Results Preview 31**

170 **Figure D-3 Synthetic Patient Record 32**

171 **Figure E-1 Calculated Fields Configuration 39**

172 **Figure E-2 Analysis Options Configuration 40**

173 **Figure E-3 Asset Details and Location 42**

174 **Figure E-4 Add/Edit Data Element 43**

175 **Figure E-5 Data Classification Email Reference Architecture 45**

176 **Figure E-6 Janusseal Schema 46**

177 **Figure E-7 Advanced Configuration Options 48**

178 **Figure E-8 Classification Profile Names and Sensitivity 49**

179 **Figure E-9 Scan Filter Configuration 51**

180 **Figure E-10 Schema Classification Settings 52**

181 **List of Tables**

182 **Table 3-1 File Types Contained in the Synthetic Dataset 6**

183 **Table 3-2 Twelve Data Types Contained in the Synthetic Dataset 6**

184 **Table 4-1 Data Classification Practices Tools and Functions 22**

185 **Table E-1 ActiveNav Products 38**

186 **Table E-2 IBM Products 41**

187 **Table E-3 Janusseal Products 44**

188 **Table E-4 Thales CipherTrust Products 47**

189 **Table E-5 Trellix Products 50**

190 **1 Overview**

191 This publication shows commercially available products performing unstructured data classification
192 practices.

193 **1.1 Challenge**

194 The project's goal is to help organizations learn how to understand its unstructured data and prepare its
195 unstructured data to leverage an organization's security and privacy controls.

196 An organization's data landscape may seem daunting and uncharted due to the sheer volume and
197 diversity of types, formats, locations, and use cases. It can encompass personally identifiable
198 information (PII) stored in databases and files on endpoints, digital conversations saved in cloud
199 environments, structured and unstructured data in data lakes and file stores, and more.

200 Since data is so vast and ubiquitous, organizations need a shared understanding of what data assets are
201 to identify and protect them.

202 The project aims at improving an organization's shared understanding of its data assets as they work to
203 meet regulatory data policies by demonstrating the use data classification tools on unstructured data.

204 The NCCoE issued an open invitation using a Federal Register Notice (FRN) [\[1\]](#) to organizations
205 interested in bringing their products and technical expertise to support and demonstrate Data
206 Classification Practices.

207 **1.2 Audience**

208 The audience of this practice guide is technical staff and leadership at medium and large organizations
209 who work with, manage or protect data. Specific audience members within these organizations may
210 include:

- 211 ▪ Chief Data Officers (CDOs) focused on digital transformation, information sharing with partners,
212 and safeguarding data
- 213 ▪ Chief Information Security Officers (CISOs)
- 214 ▪ Information security professionals focused on the data lifecycle
- 215 ▪ Data management professionals
- 216 ▪ Zero Trust Architecture (ZTA) implementers and operators
- 217 ▪ Personnel responsible for recognizing that their organization may have information in
218 unstructured data stores that needs to be protected by cryptography
- 219 ▪ Data scientists and owners who are discovering and labeling unstructured data for training
220 artificial intelligence systems.
- 221 ▪ Regulatory Compliance- support audits, meet data regulation requirements and retention policy
222 who may benefit from understanding data classification tools and are part of an organization's

- 223 ▪ Data Protection Officers whose understanding of data classification tools can support their work
224 to understand issues which relate to the protection of personal data as part of business risk and
225 liability programs.
- 226 ▪ Cost efficiency experts who wish to safeguard and protect data that are critical and delete data
227 not needed to save storage
- 228 ▪ Strategic Planners who seek to understand where data can reside like on-prem vs cloud and if
229 data can be used to support development of AI models or other innovative use cases

230 1.3 Scope

231 The scope of this publication is data classification practices, and this publication specifically
232 demonstrates discovery, identification, and labeling of unstructured data via schemas which are
233 frameworks or structures for organizing and categorizing data into different classes or categories based
234 on sensitivity, type, or business function.

235 The project does not:

- 236 • define regulatory and classification policy for specific industry sectors like USG, healthcare,
237 financial services, etc.
- 238 • demonstrate a full data management lifecycle
- 239 • perform validation of defined policies or defined schema
- 240 • demonstrate data protection practices on labelled data post data classification practices
- 241 • focus on error handling or data checking during the labeling process

242 Synthetic data was created to allow the data classification tools to work on unstructured data that
243 was free from privacy and security restrictions. The data classification tools used in this project were
244 installed and configured to look within a known network file location hosting the synthetic data.

245 1.4 Structure of This Guide

246 This NIST Cybersecurity Practice Guide provides users with the information they need to replicate the
247 discovery, identification, and labeling of data using commercially available data classification
248 technology. This guide is separated into sections, and is organized as follows:

249 [Section 2](#) provides a project overview, including the motivation for the project, the challenges in
250 implementing data classification practices, the project approach and the project’s collaborators.

251 [Section 3](#) describes the synthetic data used.

252 [Section 4](#) describes the demonstrations using steps in a workflow for each of the data classification tools
253 provided by our project collaborators.

254 [Section 5](#) documents project findings and insights.

255 Section 6 maps the tools that were used in the project to Cybersecurity Framework 2.0 subcategories.

256 **2 Project Overview**

257 **2.1 Motivation for the Project**

258 The objective of the project is to demonstrate how organizations discover, identify and label their
259 unstructured data which is a crucial component of a data management program. Data classification
260 practices enable an organization to better understand and prepare its data to be managed by data
261 protection policies that meet privacy and security regulations, to leverage Zero Trust Architectures, to
262 prioritize what should be protected by post-quantum cryptography, and to support future use of
263 Artificial Intelligence systems.

264 **2.2 Challenges in Implementing Data Classification Practices for** 265 **Unstructured Data**

266 Organizations may not know the locations of all their data, and they may not understand the sensitivity
267 or value of their unstructured data.

268 When organizations store and share their data internally and share it with other organizations, ensuring
269 the protection of the data can be challenging. Organizations that are classifying data by discovering and
270 labeling it are improving their security by understanding where their data is located, what the data's
271 sensitivity is, and what its protection needs will be.

272 Organizations that share data with other organizations lose control of that shared data once it leaves
273 their data protection boundaries. Data classification practices among data sharing organizations that
274 negotiate and use consistent labelling and tagging make it more likely that data protections can be
275 maintained at the receiving organization.

276 This publication can help organizations explore data classification practices by demonstrating:

- 277 ▪ the creation of unstructured data files from synthetic data
- 278 ▪ the use of schemas to create a logical data model that allows an organization to label its
279 unstructured data
- 280 ▪ the configuration of data classification tools to discover unstructured data that contains
281 sensitive information to be identified and labelled
- 282 ▪ the application of schemas to identify sensitive information contained within an email or
283 unstructured data using schemas
- 284 ▪ the application of labelling of sensitive information
- 285 ▪ the identifying methods to apply data classification practices to unstructured email data

286 **2.2.1 Synthetic Unstructured Data**

287 The project used synthetic data that contains realistic but fictional information. Using synthetic data
288 created from a synthetic population that statistically mirrors a real population establishes an
289 environment for experimenting with large-scale data. Using synthetic data also yields meaningful results
290 that are free from privacy and security restrictions.

291 . Identifying a source of synthetic, realistic, but not real data was an initial challenge. Creating file
292 types that contained the “not real data” for the data classification demonstrations was also necessary
293 for the project. [Section 3](#) describes the synthetic data source and using that data in over 10,000 unique
294 files.

295 **2.3 Project Approach**

296 In October 2018, the NCCoE hosted an Information Protection and Data-Centric Security Management:
297 Data Classification Workshop [\[2\]](#). The purpose of workshop was to discuss the challenges and
298 opportunities with data classification in the context of data management and information protection to
299 support various business use cases. The outcome of the workshop was used to develop goals of this
300 demonstration project. In October 2021, NIST posted a Federal Register Notice (FRN) [\[1\]](#) to solicit
301 responses from interested organizations to enter into a Cooperative Research and Development
302 Agreement (CRADA) to provide products and technical expertise to support and demonstrate security
303 platforms for the *Data Classification Practices: Facilitating Data-Centric Security Management* project.
304 Cooperative Research and Development Agreements (CRADAs) were established with selected
305 respondents to form our CRADA consortium for the project.

306 The consortium came to consensus to demonstrate unstructured data classification practices using
307 synthetic data. The project’s demonstration scenarios show how data classification techniques can
308 provide the ability to classify data using discovery, identification, and labeling techniques. Data
309 classification practices were performed on unstructured data located in a network file share location
310 and in emails.

311 **2.3.1 Project Assumptions**

312 This project is guided by the assumption that organizations would benefit from performing data
313 classification practices in a consistent and automated way that works well with the existing tools
314 because:

- 315 ▪ managing and protecting data at scale can be challenging, and data classification practices can
316 help alleviate this challenge
- 317 ▪ sharing unstructured data among two or more organizations requires coordination to identify
318 data asset types to use to describe the data that is being shared. This coordination may result in
319 an agreed upon schema, a framework or structure for organizing and categorizing data into
320 different classes or categories.
- 321 ▪ the use of an agreed upon schema between organizations will improve the likelihood that data
322 security and privacy requirements can be met within each organization even when they are not
323 using the same data classification tools.
- 324 ▪ unstructured data is in common use. Research [\[3\]](#) shows that 80-90% of organizational data is
325 unstructured and a good focus area for data classification practices.
 - 326 ○ Performing data classification practices on unstructured data can support an
327 organization’s ability to meet its privacy and data retention policies.

328 *2.3.1.1 Related Emerging Technologies*

329 This project is guided by the assumption that organizations should leverage data classification practices
330 to support:

- 331 ▪ preparing data for use by artificial intelligence systems
- 332 ▪ adoption of zero trust architectures and any security architectures which focus on securing and
333 enforcing access to data at rest and in transit through various methods, including encryption,
334 tagging and labeling, data loss prevention (DLP) strategies, and application of data rights
335 management (DRM) tools.
- 336 ▪ prioritization for migration to post-quantum cryptography

337 **3 Synthetic Data**

338 To enable demonstrations of data classification practices, synthetic unstructured data was created.

339 **3.1 Synthetic Data Characteristics**

340 The project collaborators came to consensus on use of synthetic data that has the following
341 characteristics:

- 342 • contains realistic but fictional information
- 343 • created from a synthetic population that statistically mirrors a real population
- 344 • de-identified so that it would not expose any real persons' sensitive information.
- 345 • could not be linked back to a real person's name thus avoiding privacy information
346 disclosure concerns.
- 347 • includes a range of unstructured data objects that had a known provenance or source of
348 origination
- 349 • enables standardized metadata and content
- 350 • can be used to create more than 10,000 files containing synthetic data types

351 **3.2 Synthetic Data Source**

352 SyntheticMass [4], a synthetic patient and population health data website for the state of
353 Massachusetts, was chosen as the source for creating this project's synthetic unstructured data.

354 **3.3 Unstructured Synthetic Data Files**

355 A .csv file was downloaded from SyntheticMass and used to create 25,884 files containing synthetic
356 unstructured data containing data types such as fictitious names, addresses, and patient identifiers.
357 Thirteen file types were chosen as representative of the kinds of files one might find in a file folder.
358 While scripted processes were used for creating most of the files, manual use of mail merge was also
359 used to create a small number of Word files. Visit [Appendix D Synthetic Data Creation](#) to see the process
360 steps for creating a document formatted synthetic dataset using this source.

361 Both reactive (containing sensitive data examples) and non-reactive (not containing sensitive data
 362 examples) synthetic data is included in the files to enable the data classification tools to demonstrate
 363 the ability to distinguish between sensitive and non-sensitive data types.

364 **Table 3-1 File Types Contained in the Synthetic Dataset**

| File Types | File Types |
|----------------------|---------------------|
| .doc (document) | .rtf (document) |
| .docx (document) | .csv (spreadsheet) |
| .eml (email) | .txt (document) |
| .html (web) | .wav (audio) |
| .pdf (document) | .xls (spreadsheet) |
| .png (image) | .xlsx (spreadsheet) |
| .pptx (presentation) | .zip (compact file) |

365 Synthetic data types that cover PII, Financial and Health related information were selected, as those are
 366 the dominant types found within most organizations.

367 **Table 3-2 Twelve Data Types Contained in the Synthetic Dataset**

| Data Types | Data Types |
|------------|------------------|
| Patient ID | Zip |
| Name | Birthdate |
| Address | License # |
| City | Passport # |
| County | NCCoE Customer # |
| State | NCCoE Billing # |

368 Table 3-2 shows the data types used in the 25,884 files created for the data classification practice
369 demonstrations. The 12 data types used in the demonstrations include Patient ID, Name, Address, City,
370 County, State, Zip, Birthdate, License #, Passport #, NCCoE Customer #, and NCCoE Billing #. The created
371 files may contain 0-12 of these data types. The NCCoE Customer Billing # data types were added to add
372 financially sensitive data for the data classification tools to discover. Figure 3-1 shows a specific example
373 for each data element used in the demonstrations.

SYNTHETIC DATA

Patient ID : 0a168e32-7b62-8597-0e11-296871bb764f

Name : Bryonthreeninetwo Howellninefourseven

Address : 328 Leffler Trace Unit 95

City : Chicopee

County : Hampden County

State : Massachusetts

Zip : 01013

Birthdate : 12/9/98

License # : S99997701

Passport # : X27409562X

NCCoE Customer # : NCN-741456-73442

NCCoE Billing # : NCB741456-53632

374

Figure 3-1 Specific Example of Synthetic Unstructured Data for Each Data Element

375 **4 Demonstrations**

376 The project's demonstrations show unstructured data being discovered, identified, and labeled using
377 the collaborator tools for potentially sensitive unstructured data files placed in a network file location
378 and unstructured data contained in electronic mail messages.

379 **4.1 Unstructured Data Classification Practice Demonstrations**

380 For the unstructured data classification practice demonstrations, the same 25,884 synthetic files, were
381 placed in a network file folder on host in each of the NCCoE's VMware vSphere virtualization operating
382 environments built for each data classification tool. Each product, a data classification tool, was installed
383 on a server in the same network segment as the host with the synthetic unstructured data. Each data
384 classification tool was then run to scan to discover the synthetic files in their network file location with
385 analysis done to identify the sensitive information contained in the files then apply a schema that labels
386 the data types (Patient ID, Name, Address, City, County, State, Zip, Birthdate, License #, Passport #,
387 NCCoE Customer #, and NCCoE Billing #) in each file.

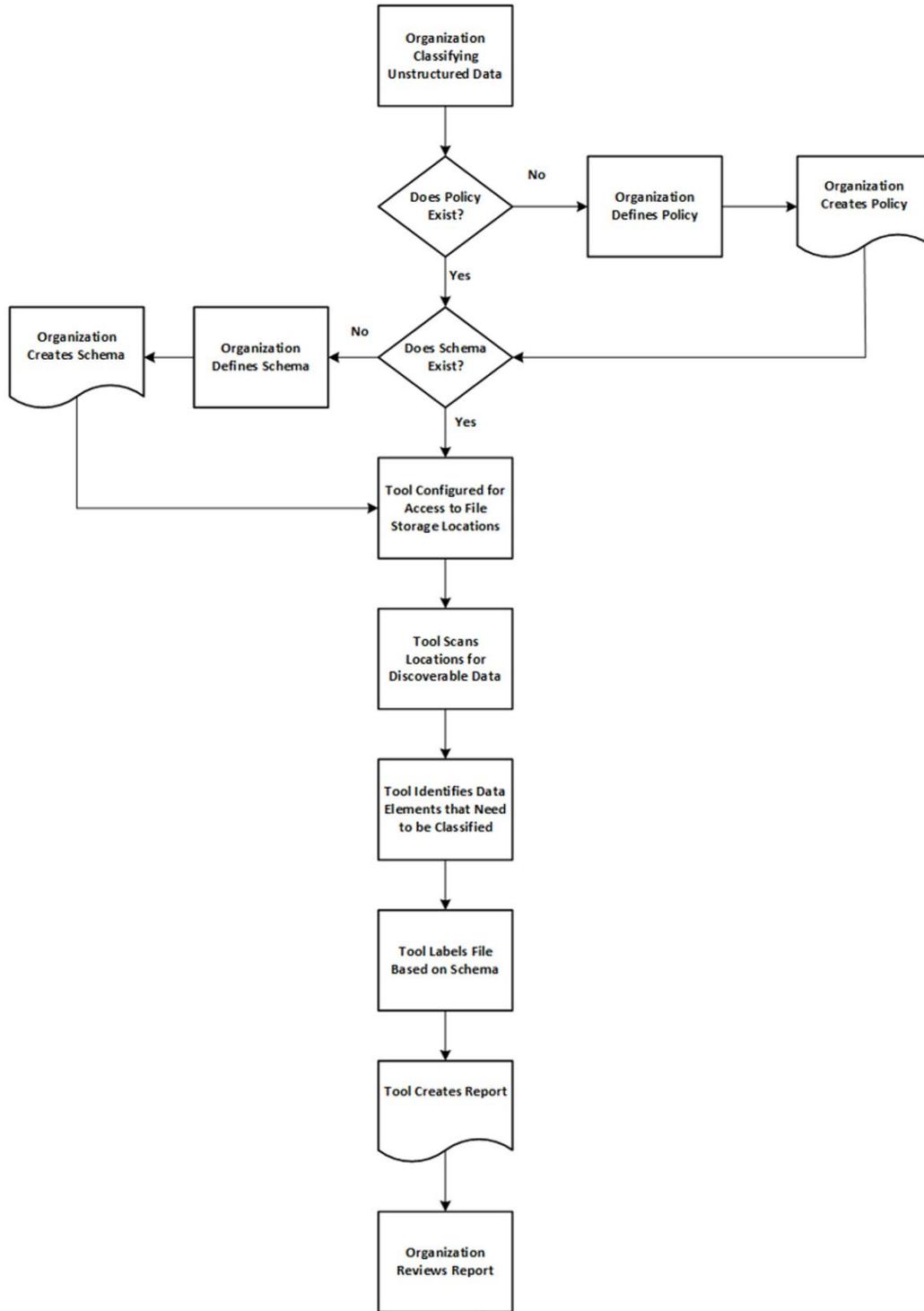
388 **4.2 Lab Demonstration Environment**

389 A separate VLAN was created for each demonstration as shown in Figure 4-1 High Level Overview of Lab
390 Environment for Each Data Classification Tool. Data Classification Tools were installed on a product
391 server. A host was created with a file storage area where the synthetic unstructured data files were
392 located. An Active Directory Server was used to provide authentication services for the Data
393 Classification Tools product servers and synthetic unstructured data file location servers. A virtual
394 desktop computer was used in each demonstration lab to configure products and provide email client
395 services for the Janusnet demonstration.

417 ▪ **Tool Labels File Based on Schema** - Labels are assigned to the data based upon its content and
418 in relation to the classification schema.

419 ▪ **Tool Creates Report** on Data Classification Practice performed on unstructured data, and
420 **Organization Reviews Report** - The tool makes the results available, and the organization
421 evaluates them.

422 The bold text above matches the text in the boxes shown below in Figure 4-2 Unstructured Data
423 Classification Practice Demonstration Workflow.



424

Figure 4-2 Unstructured Data Classification Practice Demonstration Workflow

425 4.3.1 ActiveNav Demonstration

426 ActiveNav helps organizations take control of their unstructured data. ActiveNav Discovery Center finds
427 and classifies unstructured data based on pre-defined and custom classification schemas. [Appendix E.1](#)
428 provides ActiveNav Implementation Details.

429 4.3.1.1 Schema Creation

430 For the ActiveNav demonstration scenario, a classification schema was created in ActiveNav Discovery
431 Center Workbench tool. The custom classification schema targeted the discovery of the data types
432 defined in [Table 3-2](#) including patient personal data, financial data, and health information data. The
433 classification schema was created using a rules syntax that allowed the creation of classification nodes
434 that trigger based on the content and metadata of the searched files. For example, a rule identifying all
435 files that contained one or more patient numbers was used. The classification schema could identify and
436 label files with sensitivity levels ranging from low to high based on the content they contain.

437 The classification schema was then exported from the Workbench tool and imported into ActiveNav
438 Discovery Center Project Suite application. To assign the classification schema to a file discovery scan, an
439 Index Configuration called NIST-DEMO was created that targeted the data we wanted to find in the
440 dataset, such as names, birthdates and patient billing numbers. Based on the schema, ActiveNav then
441 assigned a sensitivity level to each file that was discovered to have privacy, personal or health
442 information. With the Index Configuration created, then an Index (scan) was created that used the NIST-
443 DEMO Index Configuration as a template. The ActiveNav demonstration scenario discovered, identified,
444 and labeled data with privacy, personal credit, and health information schema categories as shown
445 below.

- 446 ▪ Privacy High
- 447 ▪ Privacy Medium
- 448 ▪ Privacy Low
- 449 ▪ Personal Credit High
- 450 ▪ Personal Credit Medium
- 451 ▪ Personal Credit Low
- 452 ▪ Health Information High
- 453 ▪ Health Information Medium
- 454 ▪ Health Information Low

455 The above schema categories could be part of an organization’s public data classification statement. For
456 example, an organization may have a policy to classify private customer data such as name and
457 birthdate with a “Privacy High” classification and classify customer metadata such as the number of
458 support tickets opened by the customer as “Privacy Medium” classification.

459 4.3.1.2 Data Discovery

460 To facilitate the scan, the synthetic dataset was copied to a network file share on a domain controller.
461 The address of the dataset was added to ActiveNav Discovery Center Project Suite as a scan location and

462 then assigned to the previously created index as a target for the data discovery. With the Index fully
463 configured, the scanning operation was initiated. The scan completed and a total of 25,884 files were
464 discovered and processed for data identification purposes.

465 *4.3.1.3 Data Identification*

466 The ActiveNav Discovery Center Index discovered and analyzed a total of 25,884 files, gathering
467 information about each one, including the file type, metadata, and content. The content of each file was
468 analyzed and compared to the classification schema to identify sensitive information. The results of this
469 process were then used to label and catalogue the files.

470 *4.3.1.4 Data Labeling / Cataloguing*

471 After the content of each discovered file was analyzed, classification labels were assigned to each
472 respective file based on the presence of data types found in the content and in relation to the
473 classification schema. For example, based upon the created schema, 10,363 files were labeled as Health
474 Information HIGH while 12,431 files were labeled with Privacy HIGH. ActiveNav Discovery Center
475 created a catalogue of all discovered files and their content. Reports can be generated to identify which
476 files contained specific pieces of data related to the classification schema.

477 *4.3.1.5 Cross-Product Integration*

478 An additional demonstration activity was performed to look at multiple organizational use cases that
479 can involve more than one product. In the additional demonstration activity, the ActiveNav product read
480 the classification tags contained in the embedded "Tags" and "Categories" custom metadata properties
481 created by Janusseau. ActiveNav used extraction rules to discover the classification tags and store them
482 as metadata associated with the file. ActiveNav then performed actions based on the classification tags
483 that Janusseau applied to the document (low, medium, high, etc.), including re-classifying the document
484 based on ActiveNav's custom rules. The classification rules in the ActiveNav product could also be used
485 to translate the classification tags found in the Janusnet documents (low, medium, high, etc.) and
486 translate them to a different classification schema such as Restricted, Confidential, and Public or
487 Internal. This demonstration showed the ability of an organization to use two products and read the
488 classification information created by one product and update the classification information using
489 another product.

490 *4.3.2 IBM Demonstration*

491 IBM Guardium Discover and Classify provides discovery and classification of structured and unstructured
492 data. [Appendix E.2](#) provides IBM Guardium Discover and Classify Implementation Details.

493 *4.3.2.1 Schema Creation*

494 Custom classification schemas were created in IBM Guardium Discover and Classify (IGDC) to represent
495 the types of data being discovered and classified. For instance, the custom data element
496 NCCOE_BILLING_NUMBER was created for this demonstration scenario to represent a patient's account
497 number. This data element was assigned a data type of "sensitive personal" and configured with a
498 custom search pattern using regex.

499 Additionally, data types can be assigned a sensitivity level such as Restricted, Confidential, and Public or
500 Internal. Once the data types have been created and configured, scans can perform discovery based on
501 the data types that have been created/configured. The IBM demonstration scenario discovered,
502 identified, and labeled data with a sensitive personal schema category across the synthetic data's 12
503 data types as shown below.

- 504 ▪ Sensitive personal schema selected across 12 data types including:
 - 505 ○ Name
 - 506 ○ Address
 - 507 ○ Birthdate

508 *4.3.2.2 Data Discovery*

509 For this demonstration scenario, the synthetic patient dataset was copied to a Server Message Block
510 (SMB) file share located on a Windows Server VM within the lab. The address of the SMB file share was
511 then added as a Data Source Catalog. The Data Source Catalog was configured with credentials to access
512 the SMB file share and assigned an analysis strategy of "baseline," which instructs the tool to perform a
513 full scan of the dataset each time. With everything configured, the scan was initiated, and IBM IGDC
514 successfully accessed the dataset in the SMB file share and processed a total of 25,884 files.

515 *4.3.2.3 Data Identification*

516 During the dataset scan, IBM IGDC processed a total of 25,884 files and analyzed them to determine
517 information about each file and the content they contained. The content of each file was compared to
518 the data types established as part of the classification schema. The results of the scan were then
519 compiled into a report detailing the number of files containing each specific data element from the
520 classification schema.

521 *4.3.2.4 Data Labeling / Cataloguing*

522 After the data in each discovered file was catalogued, a classification label was determined and assigned
523 to the appropriate files. In our schema, all discovered data types were configured as "sensitive personal"
524 data. IBM IGDC also assigned a "criticality score" to each file that helps determine the importance of the
525 data it contains. IBM IGDC also established a searchable database of all discovered files and their
526 content. Each file is assigned one or more classification labels depending on its content. For each file,
527 data element discovered in that file was provided.

528 **4.3.3 Thales Demonstration**

529 Thales CipherTrust Data Discovery and Classification (DDC) solution enables users to locate structured
530 and unstructured data. [Appendix E.3](#) provides CipherTrust DDC Implementation Details.

531 *4.3.3.1 Schema Creation*

532 In the Thales demonstration scenario, a classification profile called Synthetic Data Classification was
533 created and assigned a sensitivity level of restricted within the CipherTrust Data Discovery and
534 Classification (DDC) solution. The classification profile was assigned prebuilt and individual pieces of data
535 to search for during a scan (custom infotypes) that targeted the sensitive information found in the

536 synthetic dataset (Names, Birthdates, NCCoE_Customer_Number, Patient ID, etc.). The classification
537 profile was then assigned to a custom scan that would be used to discover the synthetic unstructured
538 dataset. The Thales demonstration scenario discovered, identified, and labeled data with a schema that
539 placed data into the following 3 categories as shown below.

- 540 ▪ Personal data
- 541 ▪ Medical
- 542 ▪ Financial

543 *4.3.3.2 Data Discovery*

544 To facilitate the scan, the synthetic patient dataset was copied to a directory on a domain controller in
545 the lab. The domain controller was then added as an available data store to be scanned by CipherTrust
546 DDC. During the creation of the custom scan, the domain controller data store was selected and
547 provided the file path where the synthetic dataset was located. With the custom scan configured, the
548 scanning operation was initiated. The scan completed and a total of 25,884 files were discovered and
549 processed.

550 *4.3.3.3 Data Identification*

551 The Thales DDC scan processed a total of 25,884 files, gathering information about each one, including
552 file type, size, modified date, and other metadata. The content of each file was analyzed and compared
553 to the classification schema to identify sensitive information. The results of this process were then used
554 to label and catalogue the files.

555 *4.3.3.4 Data Labeling / Cataloguing*

556 After the content of each discovered file was analyzed, classification labels were assigned to each
557 respective file based on its content in relation to the classification schema. For example, files containing
558 an NCCoE_Billing_Number or a passport number were labeled to show that the file contained those data
559 types. Thales DDC also created a searchable catalogue of all discovered files. Reports can also be
560 generated to identify which files contained specific pieces of data related to the classification schema.

561 *4.3.4 Trellix Demonstration*

562 Trellix Data Loss Prevention (DLP) Discover allows you to locate, classify, and protect all types of
563 corporate data. [Appendix E.4](#) provides Trellix DLP Discover implementation details.

564 *4.3.4.1 Schema Creation*

565 For the Trellix DLP Discover demonstration scenario, a classification schema was created that targeted
566 custom classifications including privacy, health info, and payment information, each with low, medium,
567 and high variations based on the number of instances found in each document. The classification
568 schema also contained out-of-the-box classifications targeting personally identifiable information and
569 protected health information. The classification policy was then assigned to a network file system (NFS)
570 file share located within the lab and automatically propagated to the DLP Discover scan. The Trellix
571 demonstration scenario discovered, identified, and labeled data with a schema that placed data into
572 privacy, health information, and Payment Card Industry (PCI) categories as shown below.

- 573 ▪ Privacy Low
- 574 ▪ Privacy Medium
- 575 ▪ Privacy High
- 576 ▪ Health Info Low
- 577 ▪ Health Info Medium
- 578 ▪ Health Info High
- 579 ▪ PCI Low
- 580 ▪ PCI Medium
- 581 ▪ PCI High

582 4.3.4.2 Data Discovery

583 The synthetic dataset was copied to the NFS file share located in the lab. The address of the dataset was
584 added to the Trellix DLP Discover scan. The scanning operation was then initiated, and the 25,884 files
585 dataset was discovered and processed.

586 4.3.4.3 Data Identification

587 The Trellix DLP Discover scan processed the synthetic file dataset, gathering information about each file,
588 including file type, size, modified date, and other metadata. The content of each file was analyzed and
589 compared to the classification schema to identify sensitive information. The results of this process were
590 then used to label and catalogue the files.

591 4.3.4.4 Data Labeling / Cataloguing

592 After the content of each discovered file was analyzed, classification labels were assigned to each
593 respective file based on their content in relation to the classification schema. For example, files
594 containing synthetic privacy information were labeled as Privacy Medium while files containing synthetic
595 health record information were labeled with Health Info Low. Trellix DLP Discover also created a
596 searchable catalogue of all discovered files. Reports can also be generated to identify which files
597 contained specific data types related to the classification schema. For example, name, address,
598 birthdate, or account numbers.

599 4.4 Electronic Mail Message Demonstration

600 Unstructured data is also commonly found in electronic mail (email) which compelled one
601 demonstration to show that unstructured data in emails can be discovered, identified, and labeled using
602 a collaborator tool.

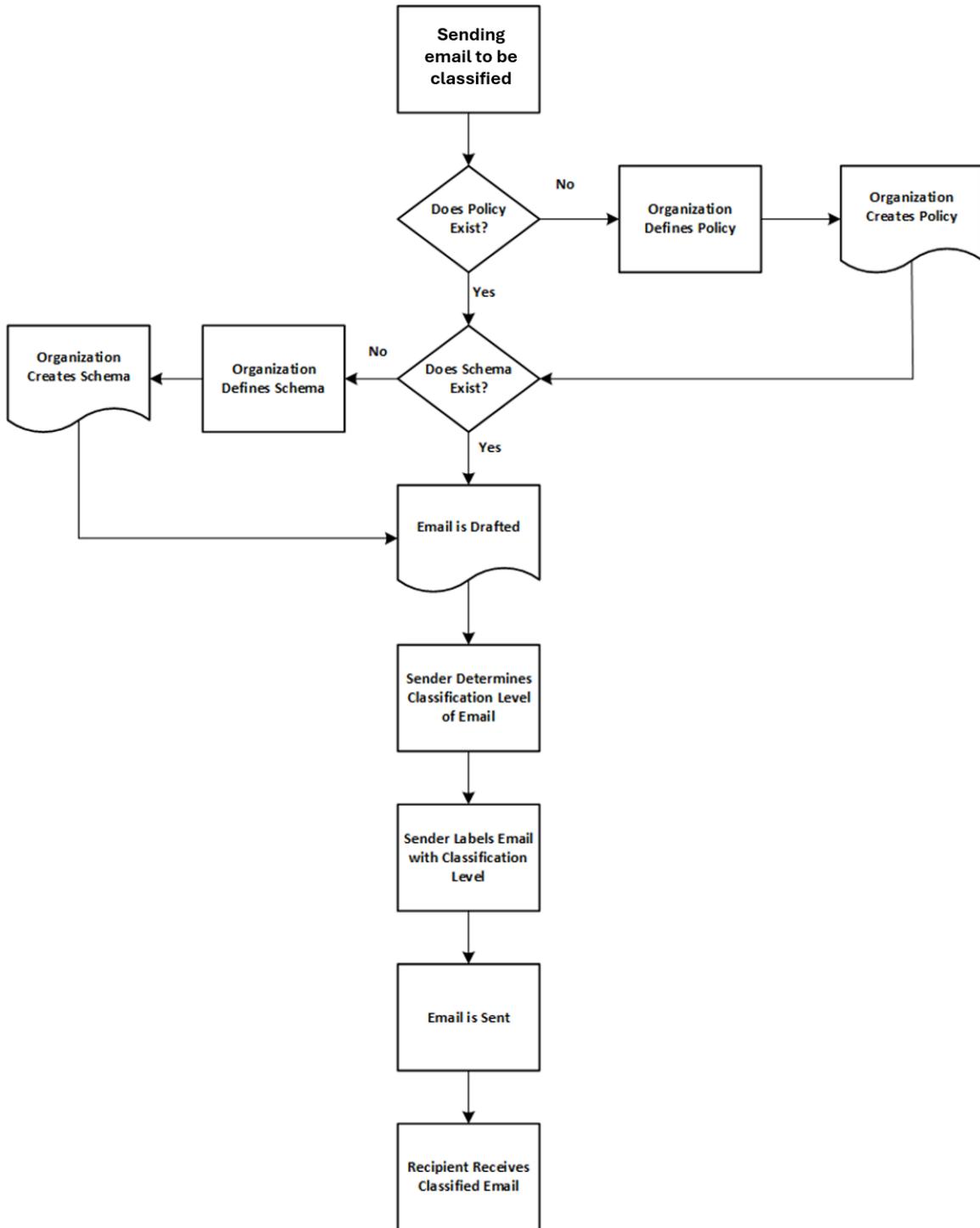
603 4.4.1 Email Practice Demonstration Workflow

604 The following describes an email data classification workflow:

- 605 ▪ **Sending Email Message to be classified** - A classified email is being sent. In the lab
606 demonstration, email is sent between 2 different accounts.

- 607 ▪ **Does Policy Exist?** Is the organization using existing policies or defining and creating their
608 classification policies? The classification policies in this publication categorize data based on the
609 data’s sensitivity. For example, an organization’s data classification policies related to the
610 sensitivity of the content of an email can require an email to be identified as “publicly
611 releasable” or “not publicly releasable”. The organization’s email data classification policies
612 define how the email should be controlled.
- 613 ▪ **Does Schema Exist?** - Is the organization using existing schema or defining and creating
614 schemas. Using the policy examples above, “Releasable”, “Low”, “Medium”, “High”, and
615 “Internal Use Only” schema hierarchy was created.
- 616 ▪ **Email is Drafted** - The sender determines the email classification level. The email’s classification
617 level is determined by the sensitivity of the message content and attachments.
- 618 ▪ **Sender Determines Classification Level of Email and Sender Labels Email with Classification**
619 **Level** - The user selects the classification level from a drop-down label that is part of the email
620 client. The labels reflect the schema that was created above. The email is labeled with
621 “Releasable”, “Low”, “Medium”, “High”, or “Internal Use Only”.
- 622 ▪ **Email is Sent and Recipient Receives the Classified Email** - Both users can see the classification
623 level that was selected for the email.

624 The workflow shown in Figure 4-3 was used to perform the demonstration steps in the lab.



625 **Figure 4-3 Email Message Demonstration Workflow**

626 **4.4.2 Janusseal Email Demonstration**

627 Janusnet products include a utility application to create separate, and different classification schemas
 628 for organizations and to perform active classification of an organization’s email messages as they were

629 sent, selecting the most appropriate level of sensitivity based on the content of the email. [Appendix E.5](#)
630 provides Janusseal Implementation Details.

631 *4.4.2.1 Use Case*

632 A common organizational use case was established with email to explore how unstructured data
633 classification needs can be improved. A single organization using a data classification product for email
634 was installed and configured. Data markings were established and communicated within a single
635 organization. Then, a multiple organization demonstration scenario using the same product in both
636 organizations was established using various schemas.

637 *4.4.2.2 Schema Creation*

638 This build demonstrated data classification practices for email. To accomplish this, two Janusnet
639 products were installed, Janusseal Schema and Janusseal for Outlook. Janusseal Schema is a utility
640 application used during the initial configuration that was used to create separate, and different
641 classification schemas for each organization to explore how organizations can interoperate with
642 different classification schemes.

643 The created schemas were then used to establish a Microsoft Group Policy Object (enterprise
644 configuration) on each organization's domain controller that propagated the Janusseal configuration to
645 users within the organization. On each user workstation, Janusseal for Outlook was installed, a plugin for
646 Microsoft Outlook, that was configured by the group policy object to allow for the user-driven
647 classification of email messages by the email sender.

648 With Janusseal for Outlook installed, users in each organization had the ability to perform active
649 classification of their email messages as they were sent, selecting the most appropriate level of
650 sensitivity based on the content of the email. The classification schemas for each organization were
651 slightly different, allowing different test scenarios to be performed by modifying the recipient of the
652 email, internal or external, or the assigned classification.

653 For the Janusnet demonstration scenario, two separate organizations were established with their own
654 domains and email servers. For each organization, Janusseal Schema was installed on a single Windows
655 workstation for use by a system administrator, and the following similar but unique classification
656 schemas were created in each fictitious company that would provide for a more realistic demonstration
657 of two organizations:

- 658 ▪ Organization A Schema (hierarchical, lowest to highest sensitivity)
 - 659 ○ Releasable
 - 660 ○ Low
 - 661 ○ Medium
 - 662 ○ High
 - 663 ○ Internal Use Only
- 664 ▪ Organizational B Schema (hierarchical, lowest to highest sensitivity)
 - 665 ○ Releasable

- 666 ○ Low
- 667 ○ Medium
- 668 ○ High
- 669 ○ Very High
- 670 ○ Not Publicly Releasable

671 Each classification schema was exported from Janusseal Schema as two separate Windows Group Policy
672 Administrative Template (ADM) files, one targeting Janusseal for Outlook and one as a Janusseal
673 product-agnostic security classification schema. With the ADM files generated, a custom Group Policy
674 Object (GPO) was created for each organization domain and then both ADM files were imported into
675 their respective domains. Each GPO was then configured to activate the Janusseal product and schemas
676 which were then applied to all Windows workstations in Organizations A and B.

677 *4.4.2.3 Data Discovery, Identification, and Labeling*

678 On each end-user Windows workstation running Microsoft Outlook Janusseal for Outlook was installed.
679 This product integrates directly into Microsoft Outlook as an Outlook add-in and enables functionality
680 based on the configuration embedded in the GPO established in the previous section. For this
681 demonstration scenario, the tool was configured to perform certain actions automatically while also
682 allowing the user to manually select a classification level to label their emails depending on content.

683 On Windows workstations in Companies 3 and 4, Microsoft Outlook was used to write emails containing
684 synthetic patient data that were then used for data discovery, identification, and labeling. Data
685 discovery and identification was performed by the user authoring the email. For the emails containing
686 the synthetic patient data, the appropriate classification level of the content was manually applied using
687 Janusseal for Outlook. The tool then added the classification markings to the top, bottom, and email
688 header metadata. For a file that is labeled using the sibling product, Janusseal Documents, when added
689 as an attachment to an email Janusseal for Outlook automatically discovers the classification of the file
690 and uses it to determine the minimum allowable classification that can be selected for the email
691 message.

692 The Janusseal products are also able to discover and write other classification metadata, such as that of
693 Microsoft's Purview labels. The Janusseal products treat its own classification metadata as authoritative,
694 if present, but can use the other metadata structures if Janusseal metadata is not present. When the
695 data has been identified, Janusnet can insert classification tags for labeling purposes into the metadata
696 of office documents and outlook emails. For Microsoft Office documents, Janusnet inserted the
697 classification tags under the "Tags" and "Categories" properties and other Office file metadata areas
698 (custom document properties).

699 **4.5 Tool Summary**

700 The name of the tools used in the Data Classification Practices demonstrations and the functions they
701 supported are shown in Table 4-1.

Table 4-1 Data Classification Practices Tools and Functions

| Tool | Data Classification Practice Function | Tool Function |
|--|--|--|
| ActiveNav Discovery Center Project Suite 4.18.0.0 | <p>Locates and classifies unstructured data, including duplicate and sensitive data discovery.</p> <p>Provides labeling to the data based upon customizable schemas that are based upon extraction rules aligned with organizational information policies.</p> | <p>Discovery Center Project Suite provides discovery, identification, and labeling of files based on the organization's policies and the data's contents and metadata. The data content determines the classification and label, which can be viewed, sorted, and filtered in configurable reports.</p> <p>Discovery Center Project Suite maintains a metadata catalog of all discovered files, and the classification(s) they were assigned.</p> <p>Discovery Center Project Suite provides incremental scans that can discover newly created files and modified content on a periodic basis.</p> <p>Discovery Center Project Suite identifies and records the presence of sensitive content within unstructured data within its metadata catalog. This supports the ability to understand and quantify risk.</p> |
| ActiveNav Discovery Center Workbench 4.18.0.0 | <p>Designs and creates the data classification policies and schemas that are implemented and used by the ActiveNav Discovery Center product.</p> | <p>Discovery Center Workbench supports the data asset selection and life cycle management with the ability to create classification schema for use in tagging files against sensitivity and other lifecycle facets.</p> |
| Janusnet Janusseal for Outlook 3.7.2.28374 | <p>Provides Microsoft Outlook users the ability to apply classification labels to email and create human-readable</p> | <p>Janusseal for Outlook can work with customizable schemas, provides classification and human-readable classification markings within the email</p> |

| Tool | Data Classification Practice Function | Tool Function |
|--|---|---|
| | classification markings within the email subject and body. | <p>subject, body, and header and the ability to block emails sent to unauthorized individuals.</p> <p>Janusseal for Outlook creates human-readable classification markings and the ability to read the classification levels the emails were assigned.</p> <p>Janusseal for Outlook provides the ability to classify email, and then reference the classifications throughout that data's lifecycle.</p> |
| Janusnet Janusseal Schema 1.4.4.24270 | Creates classification schemas that are then applied to emails and files. | Janusnet Janusseal Schema facilitates the creation of classification schemas, which define the rules determining how data will be labeled based on sensitivity and other attributes. |
| IBM Guardium Discover and Classify (IGDC) 4.0.2 | Discovery, Identification, and Labeling | <p>IGDC discovers, identifies, and labels data based on contents and metadata. Files are discovered and their content identified, and then they are labeled according to configurable policies. Results of the classification actions can be viewed, sorted, and filtered in configurable reports.</p> <p>IGDC maintains a catalog of all discovered files, their content, and the classification they were assigned.</p> <p>IGDC classifies data and provides reports that can then be used to manage data throughout its lifecycle.</p> |
| Thales CipherTrust Manager 2.17.0+12772 | Discovery, Identification, and Labeling | CipherTrust Manager provides results of the classification actions can be viewed, sorted, |

| Tool | Data Classification Practice Function | Tool Function |
|--|--|--|
| | | <p>and filtered in configurable reports.</p> <p>CipherTrust Manager maintains a catalog of all discovered files, their content, and the classification they were assigned.</p> <p>CipherTrust Manager provides reports that can then be used to manage data throughout its lifecycle.</p> |
| <p>Thales Data Discovery and Classification (DDC) 3.1.6-316</p> | <p>Discovery, Identification, and Labeling</p> | <p>Data Discovery and Classification discovers, identifies, and labels data based on contents and metadata. Files are discovered and their content identified, and then they are labeled according to configurable policies.</p> <p>Data Discovery and Classification provides incremental scans that can discover newly created files as well as existing files with new modifications.</p> |
| <p>Trellix ePolicy Orchestrator 5.10.0</p> | <p>Discovery, Identification, and Labeling</p> | <p>ePolicy Orchestrator provides results of the classification actions can be viewed, sorted, and filtered in configurable reports.</p> <p>ePolicy Orchestrator maintains a catalog of all discovered files, their content, and the classification they were assigned.</p> <p>ePolicy Orchestrator provides reports that can then be used to manage data throughout its lifecycle.</p> |

| Tool | Data Classification Practice Function | Tool Function |
|---|---|---|
| Trellix DLP Discover 11.10.500.112 | Discovery, Identification, and Labeling | DLP Discover discovers, identifies, and labels data based on contents and metadata. Files are discovered and their content identified, and then they are labeled according to configurable policies. DLP Discover provides incremental scans that can discover newly created files as well as existing files with new modifications. |

703 5 Findings and Insights

704 A listing of the general findings from the unstructured data classification practices demonstrations
705 follows:

706 **The use of synthetic data objects** was useful in our data classification practice demonstrations for the
707 following reasons and may be useful to organizations to gain experience with data classification tools:

- 708 ▪ seeing that metadata attributes in the synthetic data files were discovered by the data
709 classification tools - confidence
- 710 ▪ the creation of a variety file types with unstructured data provides the ability to test data
711 classification tool - options
- 712 ▪ having statistics about the synthetic files with data types that should be discovered by the tools -
713 verification

714 Performing a **validation of a schema** on a small subset of files before doing the whole set can help
715 avoiding multiple reruns as the classification schema is refined.

716 **Schema versioning is important.** Several versions of schemas were developed for our lab. Maintaining a
717 record of schema versions along with results of analysis on how many labels were created can help one
718 select a schema that meets your data classification needs.

719 **Appendix A List of Acronyms**

| | |
|-------|--|
| ADM | Administrative Template |
| CDO | Chief Data Officer |
| CISO | Chief Information Security Officer |
| CRADA | Cooperative Research and Development Agreement |
| DDC | Data Discovery and Classification |
| DLP | Data Loss Prevention |
| FRN | Federal Register Notice |
| GPO | Group Policy Object |
| IBM | International Business Machines Corporation |
| IGDC | IBM Guardium Discover and Classify |
| NCCoE | National Cybersecurity Center of Excellence |
| NFS | Network File System |
| NIST | National Institute of Standards and Technology |
| OCR | Optical Character Recognition |
| OMB | Office of Management and Budget |
| PCI | Payment Card Industry |
| PHI | Protected Health Information |
| PII | Personally Identifiable Information |
| SMTP | Simple Mail Transport Protocol |
| SP | Special Publication |
| ZTA | Zero Trust Architecture |

720 **Appendix B Glossary**

721 For this initial public draft, we are not prescribing terms and definitions. Please submit comments for
722 terms or words for which an explanation would enhance your understanding.

723 Appendix C References

- 724 [1] National Cybersecurity Center of Excellence (NCCoE) Data Classification Practices: Facilitating
725 Data-Centric Security Management, A Notice by the National Institute of Standards and
726 Technology on 8 October 2021. [Online]. Available:
727 [https://www.federalregister.gov/documents/2021/10/08/2021-21979/national-cybersecurity-](https://www.federalregister.gov/documents/2021/10/08/2021-21979/national-cybersecurity-center-of-excellence-nccoe-data-classification-practices-facilitating)
728 [center-of-excellence-nccoe-data-classification-practices-facilitating](https://www.federalregister.gov/documents/2021/10/08/2021-21979/national-cybersecurity-center-of-excellence-nccoe-data-classification-practices-facilitating) [Accessed 9 July 2025] .
- 729 [2] NCCoE Workshop: Information Protection and Data-Centric Security Management: Data
730 Classification Workshop , October 23, 2019 [Online]. Available:
731 [https://www.nccoe.nist.gov/get-involved/attend-events/information-protection-and-data-](https://www.nccoe.nist.gov/get-involved/attend-events/information-protection-and-data-centric-security-management-data)
732 [centric-security-management-data](https://www.nccoe.nist.gov/get-involved/attend-events/information-protection-and-data-centric-security-management-data) . [Accessed 9 July 2025]
- 733 [3] MIT Sloan stories to boost your data analytics strategy," 3 March 2021. [Online]. Available:
734 [https://mitsloan.mit.edu/ideas-made-to-matter/21-mit-sloan-stories-to-boost-your-data-](https://mitsloan.mit.edu/ideas-made-to-matter/21-mit-sloan-stories-to-boost-your-data-analytics-strategy#:~:text=Unstructured%20data%20is%20an%20area,out%20how%20to%20use%20it)
735 [analytics-](https://mitsloan.mit.edu/ideas-made-to-matter/21-mit-sloan-stories-to-boost-your-data-analytics-strategy#:~:text=Unstructured%20data%20is%20an%20area,out%20how%20to%20use%20it)
736 [strategy#:~:text=Unstructured%20data%20is%20an%20area,out%20how%20to%20use%20it](https://mitsloan.mit.edu/ideas-made-to-matter/21-mit-sloan-stories-to-boost-your-data-analytics-strategy#:~:text=Unstructured%20data%20is%20an%20area,out%20how%20to%20use%20it) .
737 [Accessed 9 July 2025]
- 738 [4] The Centers for Medicare & Medicaid Services (CMS) Alliance to Modernize Healthcare (CAMH)
739 Federally Funded Research and Development Corporation, Operated by The MITRE Corporation,
740 "SyntheticMass," [Online]. Available: <https://synthea.mitre.org/> . [Accessed 22 January 2025].

741 Appendix D Synthetic Data Creation Steps

742 The following processes were used to create the synthetic data used in the lab demonstrations. The first
 743 approach used the Mail Merge feature of Microsoft Word to create multiple Word document outputs
 744 from the synthetic data corpus while the second used a Windows PowerShell script to product a wide
 745 variety of file outputs including Word documents, Excel workbooks, PowerPoint presentations, Adobe
 746 PDFs, Rich Text Format files, HTML files, plain text files, PGN image files, WAV audio files and ZIP files.
 747 The steps of the process were performed using a Microsoft Windows based computer with Microsoft
 748 Office and the current Microsoft .NET Framework version 4 runtime installed.

749 D. 1 Download the Data

- 750 1. Visit <https://synthea.mitre.org/>
- 751 2. Select **Download Data**
- 752 3. Under **Previous Versions (2019-2021)** select **1K Sample Synthetic Patient Records, CSV**
- 753 4. Extract the files from the downloaded **.zip** file
- 754 5. Navigate to the **csv** folder and open the **patients.csv** file in Microsoft Excel. If prompted to con-
 755 vert the file's format, do not convert.

756 D. 2 Format the Data

- 757 6. (Optional) Remove **columns** that will not be part of your dataset from the patients.csv file. For
 758 example:
 - 759 a. PREFIX
 - 760 b. SUFFIX
 - 761 c. BIRTHPLACE
 - 762 d. LAT
 - 763 e. LON
 - 764 f. HEALTHCARE_EXPENSES
 - 765 g. HEALTHCARE_COVERAGE

766 D.3 Create Word Document Outputs Using Mail Merge

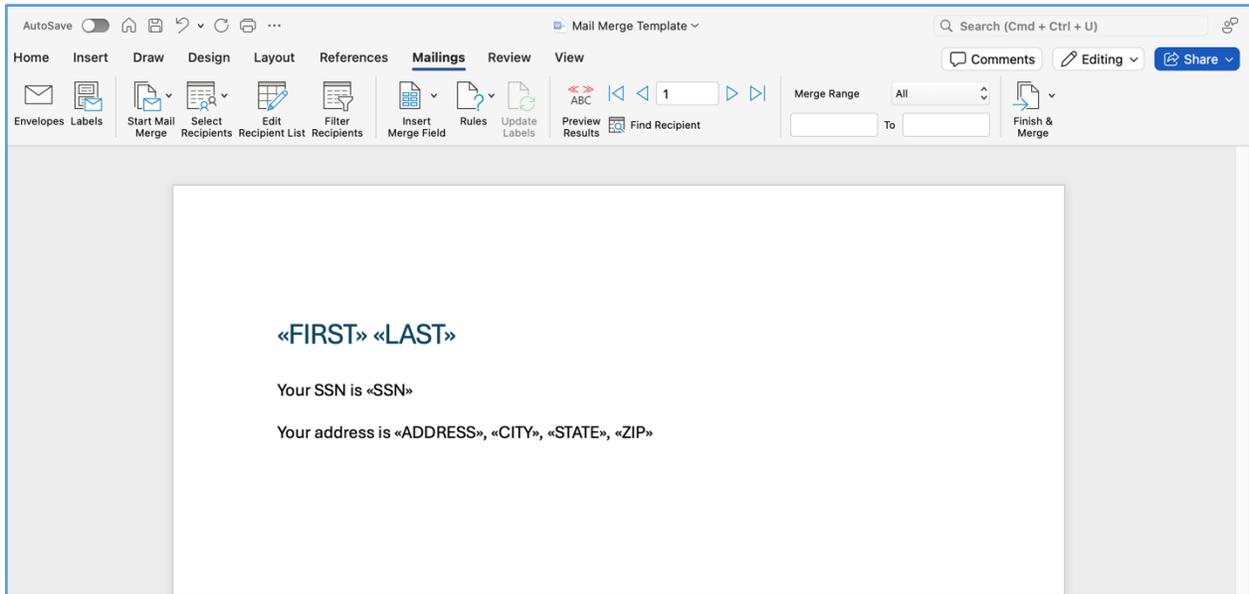
767 D.3.1 Set Up Mail Merge

- 768 7. Open a new **Microsoft Word** document
- 769 8. Click **File**, select **Save**, and name the document "**Mail Merge Template**"
- 770 9. Under the **Mailings** tab, click **Start Mail Merge** and select **Letters**
- 771 10. Click **Select Recipients** and select **Use an Existing List...**
- 772 11. Navigate to the **csv** folder containing the **synthetic patient data** from Synthea
- 773 12. Select the **patients.csv** file, click **Open**
- 774 13. If a pop-up titled "**File Conversion – patients.csv**" appears, leave the **default settings** and click
 775 **OK**

776 D.3.2 Create Mail Merge Template

- 777 14. For the first line of the document, format the text as a **heading** (found under the **Home** tab in
 778 the **Styles** section)
 - 779 a. NOTE: Formatting the first line as a header isn't necessary for the mail merge itself, but
 780 will allow us to save the merged records as individual Microsoft Word documents in a
 781 future step

- 782 15. Under the **Mailings** tab, click **Insert Merge Field** and select the data you want to include in the
783 template
- 784 a. The following field examples could be selected:
- 785 i. **First** (Name), then **insert a space** to separate first name from last name
 - 786 ii. **Last** (Name), then press **Enter** to create a carriage return
 - 787 iii. Type “Your SSN is:”, **insert a space**, then select **SSN** (Social Security Number),
788 press **Enter** to create a carriage return
 - 789 iv. Type “Your address is:”, **insert a space**, then select **Address, City, State, Zip**
- 790 b. The template will look like this:



791 **Figure D-1 Synthetic Data Template Example**

- 792 16. (OPTIONAL) Click **Preview Results** to view what the template will look like with patient data in-
793 serted
- 794 a. The results will look like this:

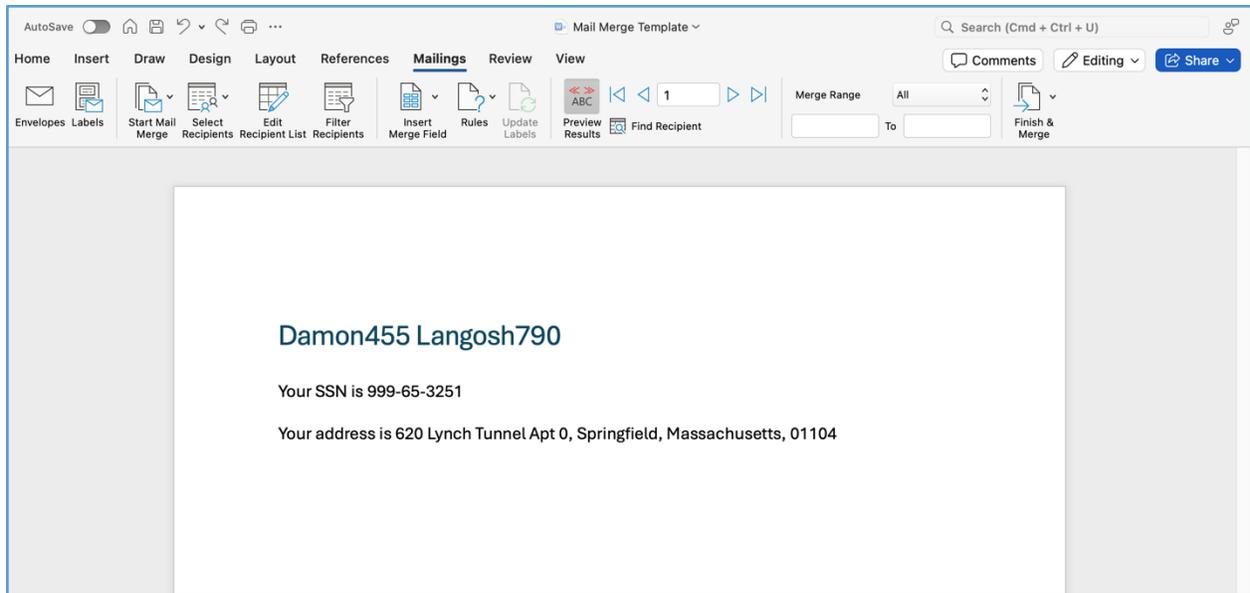


Figure D-2 Results Preview

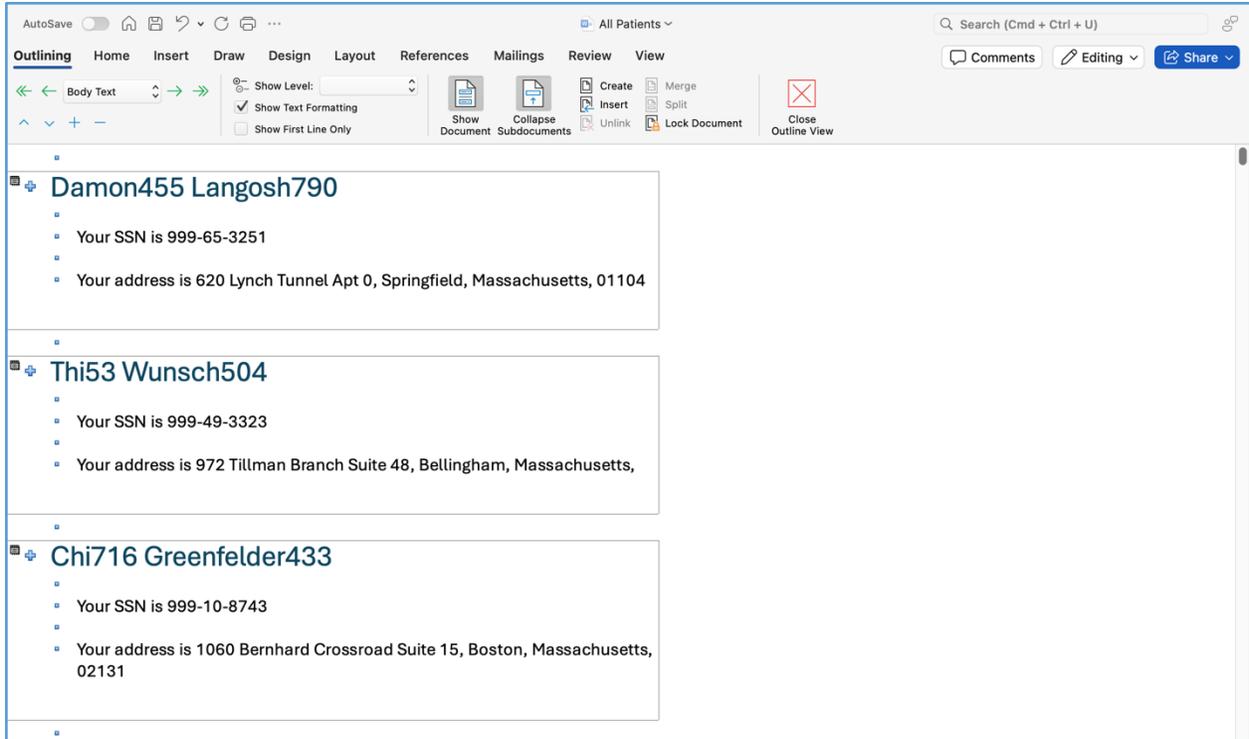
795

- 796 17. Click **Finish and Merge**, select **Edit Individual Documents**
- 797 18. Under **Merge Records**, set the value for **From:** to “1”, and for **To:** to “50”
- 798 a. This value determines which patient records are pulled from the patients.csv file. At this
- 799 point, only the first 50 records are being exported
- 800 19. Click **Ok**
- 801 20. A new Microsoft Word document will open that contains each patient record that is exported on
- 802 their own page
- 803 21. In the new document click **File**, select **Save**, and name the document “All Patients”

804

D.3.3 Create Individual Microsoft Word Documents for Each Patient

- 805 22. In the “All Patients” document, under the **View** tab, click **Outline**
- 806 23. Use **CTRL+A** to highlight all **text** in the document
- 807 24. With the **text** highlighted, click **Show Document** and then click **Create**
- 808 a. You will see data such as this:



809 **Figure D-3 Synthetic Patient Record**

- 810 25. In the “**All Patients**” document click **File** and select **Save**
- 811 a. Depending on the speed of your computer, it might take several seconds to create the
- 812 new individual documents
- 813 26. Click **Close Outline View**
- 814 27. Navigate to the folder where the “**All Patients**” document is stored
- 815 28. Find the individual Microsoft Word documents for each patient in the “**All Patients**” folder

816 **D.4 Create File Outputs Using PowerShell Script**

817 A Windows PowerShell script was developed to automate the creation of the nearly 30,000 file outputs.
 818 It is driven by the **csv** file that was obtained in the initial steps and uses various Office automation and
 819 other file creation techniques to create the outputs.

```

820 #####
821 # Author: Neville Jones, Janusnet;
822 #       neville.jones@janusnet.com
823 #
824 # No warranty
825 #
826 # Version: 0.8.0
827 # Date: 2024-09-30 15:10 +11:00
828 #####
829
830 #####
831 # Settings for running the script are here #
832 $extensions = @('doc', 'docx', 'eml', 'html', 'pdf', 'png', 'pptx', 'rtf', 'txt', 'xls', 'xlsx', 'wav', 'zip') # an array
833 $includeLabels = @( $true, $false) # whether descriptive labels (derived from CSV column headers) should be included in the generated outputs
834 $correctPersonName = $true
835 $processRowCount = -1 # the number of rows in the source CSV to process; if set to -1 then will process all rows
836 $data = Import-Csv -Path .\patients.csv -Delimiter ','
837 #####
    
```

INITIAL PUBLIC DRAFT

```
838
839 #extensions = @('txt') # for debugging - can uncomment this line to only do a single (or a subset of) extension
840
841 Add-Type -AssemblyName System.Drawing
842 Add-Type -AssemblyName System.Speech
843
844 function ReplaceDigits {
845     [OutputType([String])]
846     param ([String] $inputString)
847     $result = $inputString
848     $result = $result.Replace('0','zero')
849     $result = $result.Replace('1','one')
850     $result = $result.Replace('2','two')
851     $result = $result.Replace('3','three')
852     $result = $result.Replace('4','four')
853     $result = $result.Replace('5','five')
854     $result = $result.Replace('6','six')
855     $result = $result.Replace('7','seven')
856     $result = $result.Replace('8','eight')
857     $result = $result.Replace('9','nine')
858     return $result
859 }
860
861 $rootDir = Split-Path $MyInvocation.MyCommand.Path -Parent
862 $outputDirRoot = "output"
863 if (!(Test-Path .\$outputDirRoot)) {
864     New-Item -Path .\$outputDirRoot -ItemType "directory"
865 }
866
867 foreach($extension in $extensions) {
868     $msapp = $null
869     try {
870         Write-Host "-----"
871         foreach($includeLabel in $includeLabels) {
872
873             $outputDir = If($includeLabel) {$outputDirRoot + "\" + $extension + "\labels"} Else {$outputDirRoot + "\" + $extension + "\no labels"}
874             if (!(Test-Path .\$outputDir)) {
875                 New-Item -Path .\$outputDir -ItemType "directory"
876             }
877             Remove-Item $outputDir\*.$extension
878
879             if ($null -eq $msapp) {
880                 switch -Regex ($extension) {
881                     '^pdf$' {
882                         $_ = 'docx' # for PDF we will open an existing docx file in Word and use its ability to export PDF
883                     }
884                     '^(docx?|html|rtf)$' {
885                         $msapp = New-Object -ComObject Word.Application
886                         $msapp.Visible = $true
887                         Write-Host "---- Creating synthetic data using $($msapp.Name) ----"
888                     }
889                     '^pptx$' {
890                         $msapp = New-Object -ComObject PowerPoint.Application
891                         $msapp.Visible = 1
892                         Write-Host "---- Creating synthetic data using $($msapp.Name) ----"
893                     }
894                     '^xlsx?$' {
895                         $msapp = New-Object -ComObject Excel.Application
896                         $msapp.Visible = $true
897                         $msapp.ScreenUpdating = $false
898                         $msapp.EnableEvents = $false
899                         Write-Host "---- Creating synthetic data using $($msapp.Name) ----"
900                     }
901                 }
902             }
903
904             $i = 1
905             foreach($datum in $data) {
906                 if (($processRowCount -ne -1) -and ($i -gt $processRowCount)) {
907                     break
908                 }
909
910                 $wroteFile = $true
911                 $firstName = $datum.FIRST
912                 $lastName = $datum.LAST
```

```

913         if ($correctPersonName) {
914             $firstName = ReplaceDigits -inputString $firstName
915             $lastName = ReplaceDigits -inputString $lastName
916         }
917
918         $id = $datum.Id
919         $outputFilename = "$rootDir\$outputDir\$id.$extension"
920
921         $output = "$($datum.Disclaimer)`r`n`r`n"
922
923         $output += if($includeLabel) {"Name           : "}
924         $output += "$firstName $lastName`r`n"
925
926         $output += if($includeLabel) {"Address        : "}
927         $output += "$($datum.ADDRESS)`r`n"
928
929         $output += if($includeLabel) {"City          : "}
930         $output += "$($datum.CITY)`r`n"
931
932         $output += if($includeLabel) {"County        : "}
933         $output += "$($datum.COUNTY)`r`n"
934
935         $output += if($includeLabel) {"State         : "}
936         $output += "$($datum.STATE)`r`n"
937
938         $output += if($includeLabel) {"Zip           : "}
939         $output += "$($datum.ZIP)`r`n"
940
941         $output += if($includeLabel) {"Birthdate     : "}
942         $output += "$($datum.BIRTHDATE)`r`n"
943
944         $output += if($includeLabel) {"License #     : "}
945         $output += "$($datum.DRIVERS)`r`n"
946
947         $output += if($includeLabel) {"Passport #    : "}
948         $output += "$($datum.PASSPORT)`r`n"
949
950         $output += if($includeLabel) {"NCCoE Customer #: "}
951         $output += "$($datum.NCCoE_Customer_Number)`r`n"
952
953         $output += if($includeLabel) {"NCCoE Billing # : "}
954         $output += "$($datum.NCCoE_Billing_Number)`r`n"
955
956         switch -Regex ($extension) {
957             # simple files
958             '^eml$' {
959                 $output = "MIME-Version: 1.0`r`nContent-Type: text/plain`r`nFrom: <alice@nccoe.nist.test>`r`nTo: <bob@nccoe.nist.test>`r`nSubject:
960 test email ID $id`r`n`r`n" + $output
961             }
962             Out-File -FilePath "$outputFilename" -InputObject $output
963         }
964         '^txt$' {
965             Out-File -FilePath "$outputFilename" -InputObject $output
966         }
967         # MS Office files
968         '^(docx?|html|rtf)$' {
969             $newDoc = $msapp.Documents.Add()
970             $range = $newDoc.Paragraphs(1).Range
971             $range.Text = $output
972
973             if ($extension -eq 'docx') {
974                 $newDoc.SaveAs("$outputFilename", 16) # https://learn.microsoft.com/en-
975 us/dotnet/api/microsoft.office.interop.word.wdsaveformat?view=word-pia
976             } elseif ($extension -eq 'doc') {
977                 $newDoc.SaveAs("$outputFilename", 0)
978             } elseif ($extension -eq 'html') {
979                 $newDoc.SaveAs("$outputFilename", 10)
980             } elseif ($extension -eq 'rtf') {
981                 $newDoc.SaveAs("$outputFilename", 6)
982             }
983             $newDoc.Close()
984         }
985         '^pptx$' {
986             $newPres = $msapp.Presentations.Add()
987             $master = $newPres.SlideMaster

```

```

988     $customLayout = $null
989     foreach($cLayout in $master.CustomLayouts) {
990         if ($cLayout.Name -eq "Blank") {
991             $customLayout = $cLayout
992         }
993     }
994     $newSlide = $newPres.Slides.AddSlide(1, $customLayout)
995     $textBoxShape = $newSlide.Shapes.AddTextbox(1, 50, 50, 800, 400)
996     $textFrame = $textBoxShape.TextFrame
997     $textFrame.TextRange.Text = $output
998
999     $newPres.SaveAs("$outputFilename")
1000
1001     $newPres.Close()
1002 }
1003 '^xlsx?$' {
1004     $newWB = $msapp.Workbooks.Add()
1005     $newWS = $newWB.Worksheets[1]
1006     $newWS.Name = $($datum.Disclaimer)
1007     if ($includeLabel) {
1008         $newWS.Range("A1").Value2 = "First Name"
1009         $newWS.Range("A2").Value2 = "Last Name"
1010         $newWS.Range("A3").Value2 = "Address"
1011         $newWS.Range("A4").Value2 = "City"
1012         $newWS.Range("A5").Value2 = "County"
1013         $newWS.Range("A6").Value2 = "State"
1014         $newWS.Range("A7").Value2 = "Zip"
1015         $newWS.Range("A8").Value2 = "Birthdate"
1016         $newWS.Range("A9").Value2 = "License"
1017         $newWS.Range("A10").Value2 = "Passport"
1018         $newWS.Range("A11").Value2 = "NCCoE Customer Number"
1019         $newWS.Range("A12").Value2 = "NCCoE Billing Number"
1020     }
1021     $newWS.Range("B1").Value2 = $firstName
1022     $newWS.Range("B2").Value2 = $lastName
1023     $newWS.Range("B3").Value2 = $($datum.ADDRESS)
1024     $newWS.Range("B4").Value2 = $($datum.CITY)
1025     $newWS.Range("B5").Value2 = $($datum.COUNTY)
1026     $newWS.Range("B6").Value2 = $($datum.STATE)
1027     $newWS.Range("B7").Value2 = $($datum.ZIP)
1028     $newWS.Range("B8").Value = $($datum.BIRTHDATE)
1029     $newWS.Range("B9").Value2 = $($datum.DRIVERS)
1030     $newWS.Range("B10").Value2 = $($datum.PASSPORT)
1031     $newWS.Range("B11").Value2 = $($datum.NCCoE_Customer_Number)
1032     $newWS.Range("B12").Value2 = $($datum.NCCoE_Billing_Number)
1033
1034     if ($extension.EndsWith('x')) {
1035         $newWS.SaveAs("$outputFilename", 51) # https://learn.microsoft.com/en-us/office/vba/api/excel.xlfileformat
1036     } else {
1037         $newWS.SaveAs("$outputFilename", 56)
1038     }
1039     $result = $msapp.Workbooks.Close()
1040 }
1041 # pdf assumes that the docx extension has already been run!
1042 '^pdf$' {
1043     $sourceSubDir = If($includeLabel) {$outputDirRoot + "\docx\labels"} Else {$outputDirRoot + "\docx\no labels"}
1044
1045     $doc = $msapp.Documents.Open("$rootDir\$sourceSubDir\$id.docx")
1046     $doc.SaveAs("$outputFilename", 17)
1047     $doc.Close()
1048 }
1049 # multimedia files
1050 '^wav$' {
1051     if ($i % 100 -eq 0) { # only make audio file for every hundredth row of the CSV
1052         $speech = New-Object System.Speech.Synthesis.SpeechSynthesizer
1053         $speech.SelectVoice("Microsoft Zira Desktop")
1054         $speech.SetOutputToWaveFile($outputFilename)
1055         $speech.Speak($output)
1056         $speech.Dispose()
1057         $wroteFile = $true
1058     } else {
1059         $wroteFile = $false
1060     }
1061 }
1062 '^png$' { # https://stackoverflow.com/questions/2067920/can-i-draw-create-an-image-with-a-given-text-with-powershell

```

```

1063         if ($i % 10 -eq 0) { # only make image file for every tenth row of the CSV
1064             $bmp = New-Object System.Drawing.Bitmap 600,400
1065             $font = New-Object System.Drawing.Font Consolas,14
1066             $nistBlue = [System.Drawing.Color]::FromArgb(61,120,169)
1067             $brushBg = New-Object System.Drawing.SolidBrush $nistBlue
1068             $brushFg = [System.Drawing.Brushes]::White
1069             $graphics = [System.Drawing.Graphics]::FromImage($bmp)
1070             $graphics.FillRectangle($brushBg,0,0,$bmp.Width,$bmp.Height)
1071             $graphics.DrawString($output,$font,$brushFg,10,10)
1072             $graphics.Dispose()
1073             $bmp.Save($outputFilename)
1074             $wroteFile = $true
1075         } else {
1076             $wroteFile = $false
1077         }
1078     }
1079     # compressed files; any of these assume that the txt extension has already been run!
1080     '^zip$' {
1081         $sourceSubDir = If($includeLabel) {$outputDirRoot + "\txt\labels"} Else {$outputDirRoot + "\txt\no labels"}
1082         Compress-Archive -Path "$rootDir\$sourceSubDir\$id.txt" -CompressionLevel Fastest -DestinationPath $outputFilename
1083     }
1084
1085     default {
1086         $wroteFile = $false
1087         throw "Extension $extension not a supported output type"
1088     }
1089 }
1090
1091 if ($wroteFile) {
1092     Write-Host "Wrote record number #i to $($extension.ToUpper()) file at '$outputFilename'"
1093 }
1094
1095 $i++
1096 }
1097 }
1098 }
1099 finally {
1100     if ($null -ne $msapp) {
1101         $msapp.Quit()
1102     }
1103 }
1104 }

```

1105 D.4.2 Set up folder structure

1106 On the Windows desktop computer

- 1107 1. Open Windows Explorer
- 1108 2. Find or create a temporary folder (directory)
- 1109 3. In the temporary folder, copy the patients.csv file to the folder
- 1110 4. In the temporary folder, copy the Write-SyntheticData.ps1 PowerShell file to the folder

1111 D.4.3 Run the PowerShell script

1112 On the Windows desktop computer

- 1113 1. Open Windows Explorer
- 1114 2. Browse to the temporary folder that was created in procedure D.4.2
- 1115 3. Right-mouse click on the Write-SyntheticData.ps1 file and from the context menu select **Run with PowerShell**
- 1116 4. The PowerShell script should start and will commence creating all the necessary file outputs; it can take a significant time to run, depending on how many rows are in the source patients.csv file

- 1120 5. The output files are created in a sub-folder called “output”; within this folder are sub-folders
1121 whose name matches the file extension of the output file, so for example, under “output\xlsx”
1122 will be all the created Excel workbooks. Under each extension sub-folder there are also two sub-
1123 folders called “labels” and “no labels”; the files in the “labels” sub-folder have line delimited val-
1124 ues in the form *fieldname: synthetic data field value* while those in the “no labels” sub-folder
1125 has line delimited values in the form *synthetic data field value* (no fieldname/label information
1126 included).

1127 Appendix E Lab Implementation Details

1128 Five separate example demonstration lab builds were implemented to show unstructured data
 1129 classification needs. The example demonstrations used collaborator products that were established in
 1130 isolated virtual labs with over 60 different virtual machines in total. The virtual lab machines hosted
 1131 collaborator product installations, unstructured data locations, workstations for accessing data and that
 1132 simulated user desktop systems and included two complete enterprise-level email infrastructures. When
 1133 deploying the products into the lab environments the product installation instructions were leveraged
 1134 where possible, and information regarding the implementation details are provided for reference.

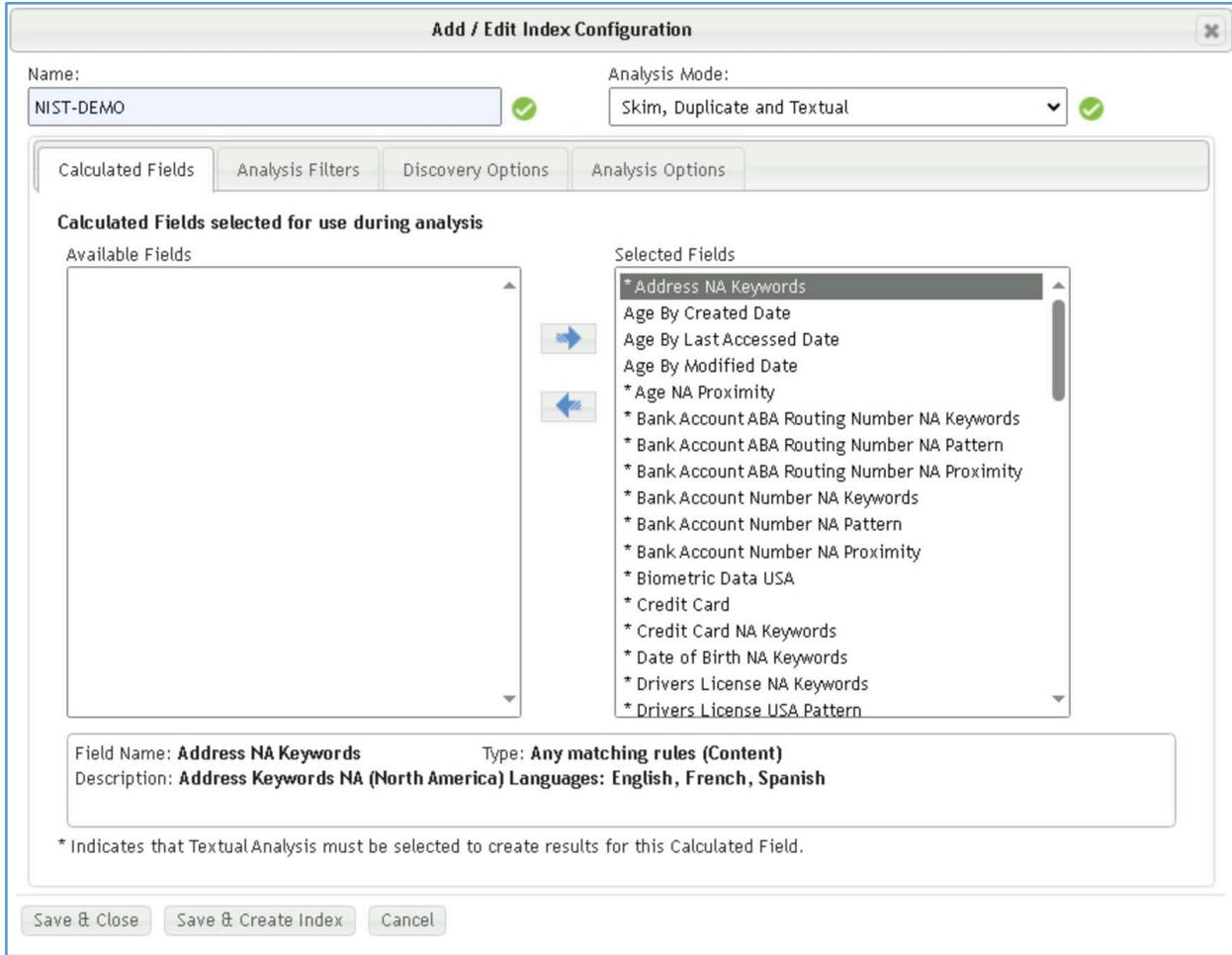
1135 E.1 ActiveNav Implementation Details

1136 ActiveNav Discovery Center finds and classifies unstructured data based on pre-defined and custom
 1137 classification schemas.

1138 **Table E-1 ActiveNav Products**

| ActiveNav Discovery Center | | Installation Notes |
|----------------------------------|---|---|
| Product Information | Product Name: Discovery Center Workbench Product Version: 4.18.0.0 Product Name: Discovery Center Project Suite Product Version: 4.18.0.0 | Installed per product instructions. |
| Configuration Information | Installation instructions available at https://support.activenav.com/hc/en-gb/articles/200738982-discovery-center-installation-guide | Synthetic data located in network file share folder. |
| Classification Schema | For information on creating a custom classification schema, see https://support.activenav.com/hc/en-gb/articles/209409665-discovery-center-starter-presentations | Schema: <ul style="list-style-type: none"> ▪ Privacy: Low, Medium, High ▪ Personal Credit: Low, Medium, High ▪ Health Information: Low, Medium, High |

1139 For this build, the on-premises version of ActiveNav Discovery Center Project Suite was used. The data
 1140 fields selected for discovery and identification were selected using the index configuration tool shown in
 1141 Figure E-1 Calculated Fields Configuration. In addition to individual field discovery and identification,
 1142 additional discovery options were configured as shown in Figure E-2 Analysis Options Configuration to
 1143 enable additional analysis options.



1144

E-1 Calculated Fields Configuration

Add / Edit Index Configuration

Name: ✓

Analysis Mode: ✓

Calculated Fields | Analysis Filters | Discovery Options | **Analysis Options**

Always re-analyze ⓘ

Enable thematic analysis ⓘ

Ignore conditional filter for duplication ⓘ

Enable content duplicate analysis ⓘ

Thematic Analysis Options

Maximum number of themes ✓ ⓘ

Maximum percentage of themes ✓ ⓘ

Number of summary sentences ✓ ⓘ

1145

Figure E-1 Analysis Options Configuration

1146 **E.2 IBM Implementation Details**

1147 IBM Guardium Discover and Classify (IGDC) provides automatic discovery and classification of structured
 1148 and unstructured data. The demonstration only exercised the discovery and classification of
 1149 unstructured data.

1150 **Table E-2 IBM Products**

| IBM Guardium Discover and Classify | | Installation Notes |
|------------------------------------|--|---|
| Product Information | Product Name: IBM Guardium Discover and Classify (IGDC) Product Version: 4.0.2 | Installed per product instructions. |
| Configuration Information | Installation instructions available at https://www.ibm.com/support/pages/ibm-guardium-discover-and-classify-documentation | Synthetic data located on network storage device. |
| Classification Schema | For information on creating a custom classification, see https://www.ibm.com/support/pages/ibm-guardium-discover-and-classify-documentation | Schema: Sensitive personal across 12 data types including: <ul style="list-style-type: none"> ▪ Name ▪ Address ▪ Birthdate |

1151 The location of the unstructured data, the data fields contained in the files, and the sensitivity levels
 1152 were configured as shown in Figure E-3 Asset Details and Location and Figure E-4 Add/Edit Data
 1153 Element.

☰ New Root Data Asset

ASSET DETAILS & LOCATION **DATA SUBJECTS**

Search for Data Subjects

Present only the mapped data elements

| Column Name | Data Element | |
|--------------------------------|--------------------------|-------------------------------|
| ✓ Disclaimer | Other ▼ | SYNTHETIC DATA |
| ✓ Id | PATIENT_ID ▼ | b9c610cd-28a6-4636-ccb6-c7... |
| ✓ BIRTHDATE | BIRTHDAY ▼ | 2/17/19 |
| ✓ NCCoE_Customer_Number | NCCOE_CUSTOMER_NUMB... ▼ | NCN-741456-72984 |
| ✓ DRIVERS | Other ▼ | S99941126 |
| ✓ PASSPORT | PASSPORT_NUMBER ▼ | X75063318X |
| ✓ FIRST | Other ▼ | Damon455 |

1154

Figure E-2 Asset Details and Location

Add/Edit Data Element ✕

Name * max 255 characters

NCCOE_CUSTOMER_NUMBER

Title * max 255 characters

NCCoE Customer Number

Sensitivity

Time value of data

Year(s)

Description max 5000 characters

Sensitive data type *

Sensitive Sensitive Personal

Data subject catalog Searchable Special categories Unique Encrypt

Classification in unstructured

Classification in unstructured with no reference data (RDA) ⓘ

Preformatters ⓘ

SpaceShrinkFormatter ✕ ▾

Persistent formatter ⓘ

SpaceShrinkFormatter ✕ ▾

1155

Figure E-3 Add/Edit Data Element

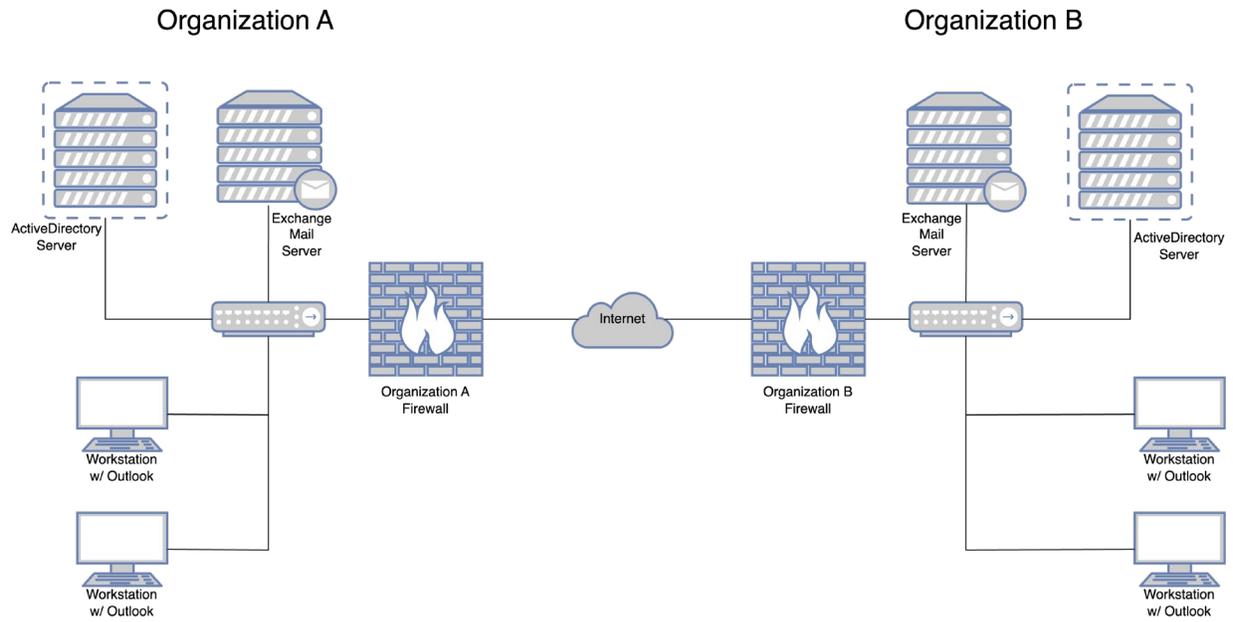
1156 **E.3 Janusseal Implementation Details**

1157 Janusseal for Outlook identifies sensitive information and enables the application of a classification label
 1158 to email. Janusseal Schema enables the creation and management of an organization's classification
 1159 schema. Janusseal was installed in a Microsoft Outlook, Windows and Exchange Server environment.

1160 **Table E-3 Janusseal Products**

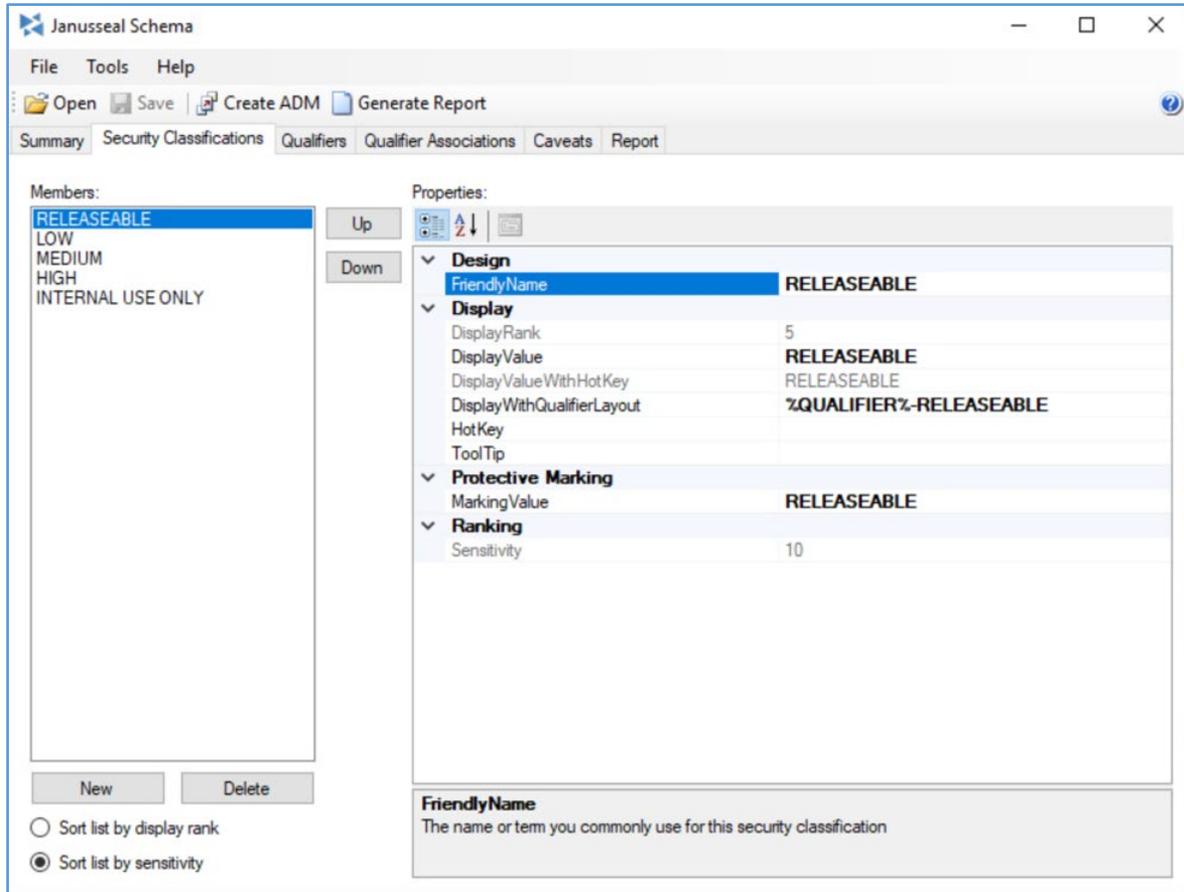
| Janusseal Products | | Installation Notes |
|----------------------------------|---|--|
| Product Information | <p>Product Name: Janusseal Schema Product Version: 1.4.4.24270</p> <p>Product Name: Janusseal for Outlook Product Version: 3.7.2.28374</p> | Installed per product instructions. |
| Configuration Information | Installation instructions available at https://www.janusnet.com/ | Products installed into two lab environments to create two separate organizations. |
| Classification Schema | <p>For information on creating a custom classification schema, see https://www.janusnet.com/</p> <p>Two slightly different schemas were implemented in the lab to create two organizations.</p> | <p>Organization A Schema:</p> <ul style="list-style-type: none"> ▪ Releasable ▪ Low ▪ Medium ▪ High ▪ Internal Use Only <p>Organization B Schema:</p> <ul style="list-style-type: none"> ▪ Releasable ▪ Low ▪ Medium ▪ High ▪ Very High ▪ Not Publicly Releasable |

1161 An architecture was developed to demonstrate the classification of unstructured data contained in
 1162 emails. The architecture enabled both single and multiple organization scenarios, where email could be
 1163 classified both within a single organization, and between two organizations as shown in Figure E-5 Data
 1164 Classification Email Reference Architecture.



1165 **Figure E-4 Data Classification Email Reference Architecture**

1166 Using the Janusseal Schema product, separate classification schemas were created for both
1167 organizations as shown in Figure E-6 Janusseal Schema. These schemas were designed to have similar,
1168 but slightly different, naming/ranking values. The schemas implemented are shown in Table E-3
1169 Janusseal Products.



1170

Figure E-5 Janusseal Schema

1171 **E.4 Thales Implementation Details**

1172 The CipherTrust platform includes Thales' data security products.

1173 **Table E-4 Thales CipherTrust Products**

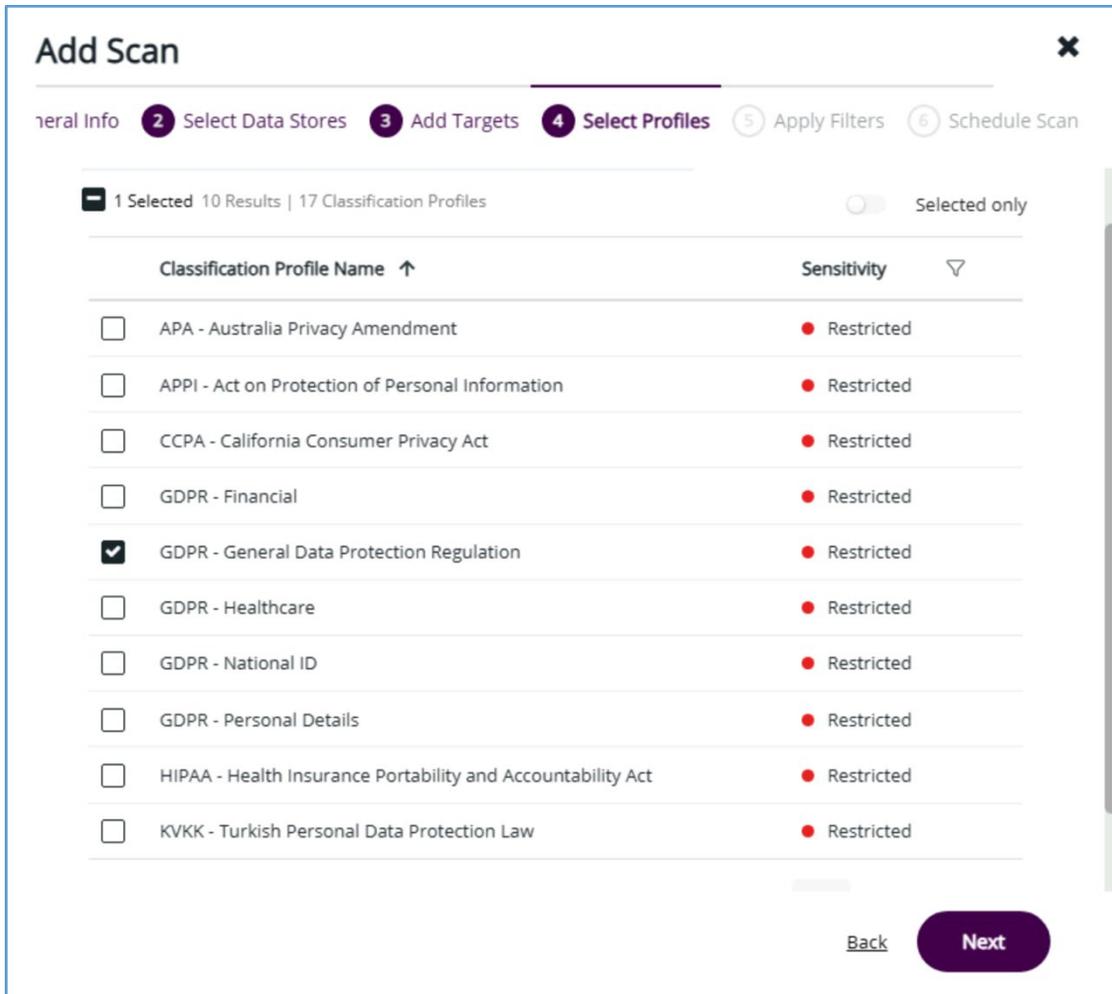
| Thales TCT CipherTrust | | Installation Notes |
|----------------------------------|--|--|
| Product Information | Product Name: Thales CipherTrust Manager Product Version: 2.17.0+12772 Product Name: Thales Data Discovery and Classification TDP Product Version: 3.1.6-316 | Installed per product instructions. |
| Configuration Information | Installation instructions available at https://www.thalestct.com/ciphertrust-data-security-platform/ https://www.thalesdocs.com/ctp/ddc/index.html | Synthetic data located on network storage device. |
| Classification Schema | For information on creating the classification schema, see https://www.thalestct.com/ciphertrust-data-security-platform/ | Schema: <ul style="list-style-type: none"> ▪ Personal data ▪ Medical ▪ Financial |

1174 The file types to be scanned by the Thales TCT CipherTrust tools, their classification profiles and
 1175 sensitivity were selected as shown in Figure E-7 Advanced Configuration Options and Figure E-8
 1176 Classification Profile Names and Sensitivity.

The screenshot shows a dialog box titled "Add Scan" with a close button (X) in the top right corner. Below the title is a progress indicator with six steps: 1. General Info (selected), 2. Select Data Stores, 3. Add Targets, 4. Select Profiles, 5. Apply Filters, and 6. Schedule. The main content area includes a text input field for an "Optional description of up to 250 characters". Below this is a section titled "Advanced Configuration" with a dropdown arrow. Under "Advanced Configuration", there are several settings: "Scan Priority" with radio buttons for "Low" and "Normal" (selected); "Content supported:" with checkboxes for "OCR", "Voice", and "EBCDIC" (all checked); "Trace logs:" with a checkbox for "Logs" (checked); "Memory Usage Limit (MB):" with an input field containing "1024"; "Throughput (MBps):" with an input field containing "0"; and "Amount of Data Object Volume:" with radio buttons for "Low - maximum info", "Medium - core info", "High - minimal info", and "High - no saved info" (selected). At the bottom left is a "Restore Defaults" button, and at the bottom right are "Cancel" and "Next" buttons.

1177

Figure E-6 Advanced Configuration Options



1178

Figure E-7 Classification Profile Names and Sensitivity

1179 **E.5 Trellix Implementation Details**

1180 Trellix Data Loss Prevention (DLP) supports agent-based and agentless discovery and classification of
 1181 data-at-rest on workstation endpoints and network resources.

1182 **Table E-5 Trellix Products**

| Trellix Data Loss Prevention | | Installation Notes |
|----------------------------------|---|---|
| Product Information | <p>Product Name: Data Loss Prevention Discover Server Product Version: 11.10.500.112</p> <p>Product Name: Trellix ePolicy Orchestrator Product Version: 5.10.0</p> | Installed per product instructions. |
| Configuration Information | <p>Trellix Data Loss Prevention Discover Server installation instructions can be found at https://docs.trellix.com/bundle</p> <p>Trellix ePolicy Orchestrator installation instructions can be found at https://docs.trellix.com/bundle/epolicy-orchestrator-landing/page/UUID-52b91793-68f9-a8ac-141c-104824763f9a.html</p> | Synthetic data located on network storage device. |
| Classification Schema | For information on creating a custom classification schema, see https://docs.trellix.com/bundle/data-loss-prevention-11.1.x-product-guide/page/GUID-F7DAE857-DC7C-44DA-AEB4-BFFB1280FD0B.html | <p>Schema:</p> <ul style="list-style-type: none"> ▪ Privacy Low ▪ Privacy Medium ▪ Privacy High ▪ Health Info Low ▪ Health Info Medium ▪ Health Info High ▪ PCI Low ▪ PCI Medium ▪ PCI High |

1183 Scan filters were applied to identify individual data fields in the synthetic unstructured data as shown in
 1184 Figure 0-9 Scan Filter Configuration. Sensitivity of the scanned data was configured using a filter as
 1185 shown in Figure 0-10 Schema Classification Settings.

Choose from existing values.

Check one or more items from the list:

Filter items: Show selected items only. Include Built-in items

| Name | Threshold | Actions |
|---|-----------|----------------------|
| <input type="checkbox"/> Mexico Voter Card Number [built-in] | 1 | View |
| <input type="checkbox"/> Mexico banking standard (CLABE) [built-in] | 1 | View |
| <input type="checkbox"/> Microsoft Windows Product Key [built-in] | 1 | View |
| <input type="checkbox"/> Monero (XMR) Address [built-in] | 1 | View |
| <input type="checkbox"/> Multiple Common PCI Cards [built-in] | 1 | View |
| <input type="checkbox"/> NCCoE Billing Number | 1 | Edit |
| <input checked="" type="checkbox"/> NCCoE Customer Number | 1 | Edit |
| <input type="checkbox"/> National Drug Code (NDC) [built-in] | 1 | View |
| <input type="checkbox"/> National Provider Identifier (HIPAA) [built-in] | 1 | View |
| <input type="checkbox"/> New Zealand IRD Number [built-in] | 1 | View |
| <input type="checkbox"/> New Zealand Ministries of Health Index Number [built-in] | 1 | View |

Count multiple occurrences of each match string.
 Count each match string only one time.

1186

Figure E-8 Scan Filter Configuration

Choose classifications

Select one or more classifications from the list. ?

Filter items: **GO** Show selected items only. Include Built-in items

| <input type="checkbox"/> | Name |
|-------------------------------------|--------------------|
| <input type="checkbox"/> | Partner Data High |
| <input type="checkbox"/> | Partner Data Low |
| <input checked="" type="checkbox"/> | Privacy Low |
| <input checked="" type="checkbox"/> | Privacy Medium |
| <input checked="" type="checkbox"/> | Privacy High |
| <input checked="" type="checkbox"/> | Health Info Low |
| <input checked="" type="checkbox"/> | Health Info Medium |
| <input checked="" type="checkbox"/> | Health Info High |
| <input checked="" type="checkbox"/> | PCI Medium |
| <input checked="" type="checkbox"/> | PCI High |
| <input checked="" type="checkbox"/> | PCI Low |
| <input type="checkbox"/> | NCCAF High |

New Classification **Add**

OK **Cancel**

1187

Figure E-9 Schema Classification Settings