



1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19

20
21

**NIST Internal Report
NIST IR 8432 ipd**

Cybersecurity of Genomic Data

Initial Public Draft

Ron Pulivarti
Natalia Martin
Fred Byers
Justin Wagner
Justin Zook
Samantha Maragh
Kevin Wilson
Martin Wojtyniak
Brett Kreider
Ann-Marie France
Sallie Edwards
Tommy Morris
Jared Sheldon
Scott Ross
Phillip Whitlow

This publication is available free of charge from:
<https://doi.org/10.6028/NIST.IR.8432.ipd>

22
23
24

NIST Internal Report NIST IR 8432 ipd Cybersecurity of Genomic Data

25 Initial Public Draft

26 Ron Pulivarti
27 Natalia Martin
28 Fred Byers
29 *National Cybersecurity*
30 *Center of Excellence*
31 *Information Technology*
32 *Laboratory*

33
34 Justin Wagner
35 Justin Zook
36 Samantha Maragh
37 *Material Measurement*
38 *Laboratory*

Kevin Wilson
Martin Wojtyniak
Brett Kreider
Ann-Marie France
Sallie Edwards
MITRE

Tommy Morris
Jared Shelton
University of Alabama in Huntsville

Scott Ross
Phillip Whitlow
HudsonAlpha

41 This publication is available free of charge from:
42 <https://doi.org/10.6028/NIST.IR.8432.ipd>

43 March 2023



44 U.S. Department of Commerce
45 *Gina M. Raimondo, Secretary*

46 National Institute of Standards and Technology
47 *Laurie E. Locascio, NIST Director and Under Secretary of Commerce for Standards and Technology*

48 Certain commercial entities, equipment, or materials may be identified in this document in order to describe an
49 experimental procedure or concept adequately. Such identification is not intended to imply recommendation or
50 endorsement by the National Institute of Standards and Technology (NIST), nor is it intended to imply that the
51 entities, materials, or equipment are necessarily the best available for the purpose.

52 There may be references in this publication to other publications currently under development by NIST in
53 accordance with its assigned statutory responsibilities. The information in this publication, including concepts and
54 methodologies, may be used by federal agencies even before the completion of such companion publications. Thus,
55 until each publication is completed, current requirements, guidelines, and procedures, where they exist, remain
56 operative. For planning and transition purposes, federal agencies may wish to closely follow the development of
57 these new publications by NIST.

58 Organizations are encouraged to review all draft publications during public comment periods and provide feedback
59 to NIST. Many NIST cybersecurity publications, other than the ones noted above, are available at
60 <https://csrc.nist.gov/publications>.

61 **NIST Technical Series Policies**

62 [Copyright, Use, and Licensing Statements](#)

63 [NIST Technical Series Publication Identifier Syntax](#)

64 **How to Cite this NIST Technical Series Publication:**

65 Pulivarti R, Martin N, Byers F, Wagner J, Maragh S, Wilson K, Wojtyniak M, Kreider B, Frances A, Edwards S,
66 Morris T, Sheldon J, Ross S, Whitlow P (2023) Cybersecurity of Genomic Data. (National Institute of Standards and
67 Technology, Gaithersburg, MD), NIST Interagency or Internal Report (IR) NIST IR 8432 ipd.
68 <https://doi.org/10.6028//NIST.IR.8432.ipd>

69 **Author ORCID iDs [will be updated in the final publication]**

70 Ron Pulivarti: 0000-0002-8330-3474

71 Natalia Martin: 0000-0000-0000-0000

72 Fred Byers: 0000-0000-0000-0000

73 Justin Wagner: 0000-0000-0000-0000

74 Justin Zook: 0000-0003-2309-8402

75 Samantha Maragh: 0000-0003-2564-9589

76 **Public Comment Period**

77 March 3, 2023 - April 3, 2023

78 **Submit Comments**

79 genomic_cybersecurity_nccoe@nist.gov

80 National Institute of Standards and Technology

81 Attn: Applied Cybersecurity Division, Information Technology Laboratory

82 100 Bureau Drive (Mail Stop 2002) Gaithersburg, MD 20899-8930

83 **All comments are subject to release under the Freedom of Information Act (FOIA).**

84 **Abstract**

85 Genomic data has enabled the rapid growth of the U.S. bioeconomy and is valuable to the
86 individual, industry, and government because it has multiple intrinsic properties that in
87 combination make it different from other types of high value data which possess only a subset of
88 these properties. The characteristics of genomic data compared to other high value datasets raises
89 some correspondingly unique cybersecurity and privacy challenges that are inadequately
90 addressed with current policies, guidance documents, and technical controls.

91 This report describes current practices in risk management, cybersecurity, and privacy
92 management for protecting genomic data along with relevant challenges and concerns. Gaps in
93 protection practices across the lifecycle were identified concerning genomic data generation, safe
94 and responsible sharing of the genomic data, monitoring the systems processing genomic data,
95 lack of specific guidance documents addressing the unique needs of genomic data processors,
96 and regulatory/policy gaps with respect to national security and privacy threats in the collection,
97 storage, sharing, and aggregation of human genomic data.

98 The report proposes a set of solution ideas that address real-life use cases occurring at various
99 stages of the genomic data lifecycle along with candidate mitigation strategies and the expected
100 benefits of the solutions. Additionally, areas needing regulatory/policy enactment or further
101 research are highlighted.

102 **Keywords**

103 Cyberbiosecurity; cybersecurity; genomic data; human genome; genomics; privacy.
104

105 **Reports on Computer Systems Technology**

106 The Information Technology Laboratory (ITL) at the National Institute of Standards and
107 Technology (NIST) promotes the U.S. economy and public welfare by providing technical
108 leadership for the Nation's measurement and standards infrastructure. ITL develops tests, test
109 methods, reference data, proof of concept implementations, and technical analyses to advance
110 the development and productive use of information technology. ITL's responsibilities include the
111 development of management, administrative, technical, and physical standards and guidelines for
112 the cost-effective security and privacy of other than national security-related information in
113 federal information systems.

114 **Additional Information**

115 For additional information on this NIST Internal Report and the NCCoE Genomics
116 Cybersecurity project, visit [our project page](#). Information on other efforts at [NIST](#), the
117 [Information Technology Laboratory](#) and [NCCoE](#) is also available.

118 **Call for Patent Claims**

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

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

- 127 a) assurance in the form of a general disclaimer to the effect that such party does not hold
128 and does not currently intend holding any essential patent claim(s); or
- 129 b) assurance that a license to such essential patent claim(s) will be made available to
130 applicants desiring to utilize the license for the purpose of complying with the guidance
131 or requirements in this ITL draft publication either:
 - 132 i. under reasonable terms and conditions that are demonstrably free of any unfair
133 discrimination; or
 - 134 ii. without compensation and under reasonable terms and conditions that are
135 demonstrably free of any unfair discrimination.

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

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

143 Such statements should be addressed to: genomic_cybersecurity_nccoe@nist.gov

144	Table of Contents	
145	1. Introduction	1
146	1.1. Cybersecurity and Privacy Concerns	1
147	1.2. Document Scope and Goals	2
148	2. Background	3
149	2.1. Genomic Information Lifecycle	3
150	2.2. Next Generation Sequencing	5
151	2.3. Variant Calling	5
152	2.4. Genome Editing	5
153	2.5. Direct-to-Consumer Testing	6
154	2.6. The Characteristics of Genomic Data	6
155	2.7. Balance Between Benefits and Risks for Uses of Genomic Information	7
156	3. Challenges and Concerns Associated with Handling Genomic Information	9
157	3.1. Potential National Security Concerns	9
158	3.2. Privacy Challenges	9
159	3.3. Discrimination and Reputational Concerns	10
160	3.4. Economic Concerns	10
161	3.5. Health Outcome Concerns	11
162	3.6. Other Potential Future Concerns	11
163	3.7. Summary of Challenges and Concerns with Genomic Data	11
164	4. Current State of Practices	12
165	4.1. Risk Management Practices	12
166	4.1.1. U.S. Government Resources	12
167	4.1.2. U.S. Government Initiatives	14
168	4.1.3. International Resources and Regulations	15
169	4.2. Cybersecurity Best Practices	15
170	4.2.1. U.S. Government Resources	15
171	4.2.2. International Resources	16
172	4.2.3. Industry Resources	16
173	4.3. Privacy Best Practices	17
174	4.3.1. U.S. Government Resources	17
175	4.3.2. Industry Resources	18
176	4.3.3. International Resources	18
177	4.4. Summary of Gaps in the Protection of Genomic Data	18
178	4.4.1. Guidance Gaps	18

179	4.4.2. Technical Solution Gaps	19
180	4.4.3. Policy/Regulatory Landscape	20
181	5. Available Solutions to Address Current Needs	21
182	5.1. NIST Cybersecurity Framework Profile for Processing of Genomic Data	21
183	5.1.1. Use Case Description.....	21
184	5.1.2. Solution Idea	21
185	5.1.3. Expected Benefits	22
186	5.2. NIST Privacy Framework Profile for Processing Genomic Data	22
187	5.2.1. Use Case Description.....	22
188	5.2.2. Solution Idea	23
189	5.2.3. Expected Benefits	23
190	5.3. Automatic Network Micro-Segmentation of Sequencers with MUD	23
191	5.3.1. Use Case Description.....	23
192	5.3.2. Solution Idea	23
193	5.3.3. Expected Benefits	24
194	5.4. Security Guidelines for Data Analysis Pipelines	24
195	5.4.1. Use Case Description.....	25
196	5.4.2. Solution Idea	25
197	5.4.3. Expected Benefits	25
198	5.5. Demonstration Project for Genomic Data Risk Management.....	25
199	5.5.1. Use Case Description.....	26
200	5.5.2. Solution Idea	26
201	5.5.3. Expected Benefits	26
202	5.6. Demonstration Project for Analysis of Genomic Data Using Privacy Preserving and	
203	Enhancing Technologies	27
204	5.6.1. Use Case Description.....	28
205	5.6.2. Solution Idea	28
206	5.6.3. Expected Benefits	29
207	6. Areas for Further Research.....	30
208	7. Conclusion	31
209	8. References.....	32
210	Appendix A. List of Symbols, Abbreviations, and Acronyms.....	39

211	List of Tables	
212	Table 1. NIST Risk Management Framework Discussion for Genomic Data.....	13
213	List of Figures	
214	Figure 1. Genomic Data Lifecycle; Naveed et.al. [6].....	3
215	Figure 2. Characteristics of DNA (Adapted from Figure 1 of Naveed et al. [6]).....	7
216	Figure 3. Notional MUD architecture for a sequencer.....	24

217 **Acknowledgments**

218 The authors would like to thank the speakers, panelists, and contributors to the NCCoE Virtual
219 Workshop on the Cybersecurity of Genomic Data held Wednesday, January 26, 2022, and the
220 NCCoE Virtual Workshop on Exploring Solutions for the Cybersecurity of Genomic Data held
221 Wednesday, May 18, 2022, and Thursday, May 19, 2022, for their inputs, which helped inform
222 the contents of this paper.

223 **1. Introduction**

224 The world has entered an era of accelerated biological innovation built primarily upon the many
225 uses of genomic data that include vaccine development and manufacturing, pharmaceutical
226 development and manufacturing, disease diagnosis, agricultural innovations that enable
227 increased food production, biofuel development, basic and translational scientific research,
228 consumer testing, genealogy, and law enforcement, among others. More uses continue to be
229 discovered. Genetic sequencing technology has advanced such that sequencing entire genomes is
230 feasible and affordable. Whole or partial genome sequences for many microbial, plant, and
231 animal species reside in open access, controlled access, or private databases within the National
232 Institutes of Health (NIH), Federal Bureau of Investigation (FBI), and direct-to-consumer (DTC)
233 genetic testing providers, to name a few. As this era unfolds there is a new awareness of risks to
234 U.S. national security, its economy, its biotechnology industry, and its citizens due to
235 cybersecurity attacks targeting genomic data as highlighted in the Executive Order on Advancing
236 Biotechnology and Biomanufacturing Innovation for a Sustainable, Safe, and Secure American
237 Bioeconomy [1]. Additionally, human genetic information requires complying with policies,
238 laws, and ethics surrounding privacy. Nevertheless, the inherent value of some genomic data lies
239 in the ability to share information with the broader community, creating the need to balance
240 access restrictions with data sharing capabilities.

241 **1.1. Cybersecurity and Privacy Concerns¹**

242 Cyber attacks targeted at genomic data include attacks against the confidentiality of the data, its
243 integrity, and its availability. Cyber attacks against the confidentiality of the data can threaten
244 our economy through theft of the intellectual property owned by the U.S. biotechnology industry,
245 allowing competitors to gain an unfair economic advantage by accessing U.S. held genomic data.
246 Attacks against the integrity of the data can disrupt biopharmaceutical output, agricultural food
247 production, and bio-manufacturing activity. Attacks against the availability of the data include
248 encrypting for ransom, deletion of data, and disabling critical automated equipment used in
249 research, development, and manufacturing. The potential harms of cyber attacks on genomic data
250 threaten our national security as well, including enabling the development of biological weapons
251 and the surveillance, oppression, and extortion of our citizens, military, and intelligence
252 personnel based on their genomic data.

253 Cyber attacks targeted at genomic data can also harm individuals by enabling blackmail,
254 discrimination based on disease risk, and privacy loss from the revealing of hidden consanguinity
255 or phenotypes including health, emotional stability, mental capacity, appearance, and physical
256 abilities. In addition to the privacy risks that can arise because of a cyber attack, privacy risks
257 unrelated to cybersecurity can arise when processing genomic data. These risks can arise when

¹ Cybersecurity and privacy objectives provide a useful construct for describing concerns. The cybersecurity objectives are confidentiality, integrity, and availability. The definitions of each originate from 44 U.S.C., Sec. 3542 and they are used throughout multiple NIST publications. The privacy engineering objectives are predictability, manageability, and disassociability. They were initially described in NISTIR 8062 and are also used throughout multiple NIST publications. Definitions for all six terms, as well as a sample of documents in which they are discussed, appear in the NIST Glossary (available at: <https://csrc.nist.gov/glossary>).

258 there is insufficient predictability, manageability, and disassociability in the genomic data
259 processing. Insufficient predictability in data processing can result in privacy problems if
260 individuals are surprised by what is happening with their genomic data. Insufficient
261 manageability in data processing can arise when the capabilities are not in place to allow for
262 appropriately granular administration of genomic data, for example, individuals may need to be
263 able to have some or all their genomic data deleted from a dataset. Permitting access to raw
264 genomic data, instead of using appropriate privacy-enhancing technologies to extract only the
265 necessary insights (without revealing the raw data), introduces privacy risks from insufficient
266 disassociability in data processing. Each of these areas of privacy risks can disrupt the ability to
267 realize the benefits of processing genomic data [2].

268 NIST is exploring genomic data uses to better understand common and pressing cybersecurity
269 and privacy concerns specific to this data to identify and provide security and privacy practice
270 guidance to help protect it. To inform this effort, NIST, through its National Cybersecurity
271 Center of Excellence (NCCoE), conducted two public workshops [3][4]—one concentrating on
272 the challenges already faced or anticipated by the community, followed by another workshop
273 focusing on solutions to help address those challenges. These workshops along with additional
274 research provided the basis for the information presented in this paper.

275 **1.2. Document Scope and Goals**

276 This paper identifies characteristics of genomic data compared to other types of high-value data
277 and provides an introductory overview of cybersecurity and privacy risk management resources
278 and classifications of the risks during the genomic data lifecycle. It also identifies the most
279 common challenges to securing genomic data, the current state-of-the-art cybersecurity and
280 privacy practices for genomic data, and gaps associated with these practices. Finally, a set of use
281 cases are presented that simulate real life challenges with candidate mitigation strategies to
282 address each challenge, along with the expected benefits provided by the proposed solutions.

283 2. Background

284 A genome contains hereditary material comprised of nucleic acids, mostly in the form of
285 deoxyribonucleic acid (DNA). Genomes contain the full set of instructions to form an organism
286 and are largely unchanged from conception to death. Instructions are encoded in the sequence of
287 the four nucleotide subunits (also called bases) that comprise the backbone of the DNA or
288 Ribonucleic acid (RNA) molecule. In DNA, these are adenine (A), cytosine (C), guanine (G),
289 and thymine (T). A segment of bases containing the instructions for making a product, such as a
290 protein or an RNA molecule, is called a gene. Variations in the number and kinds of genes, as
291 well as differences across genes, underpin the diversity of life on earth.

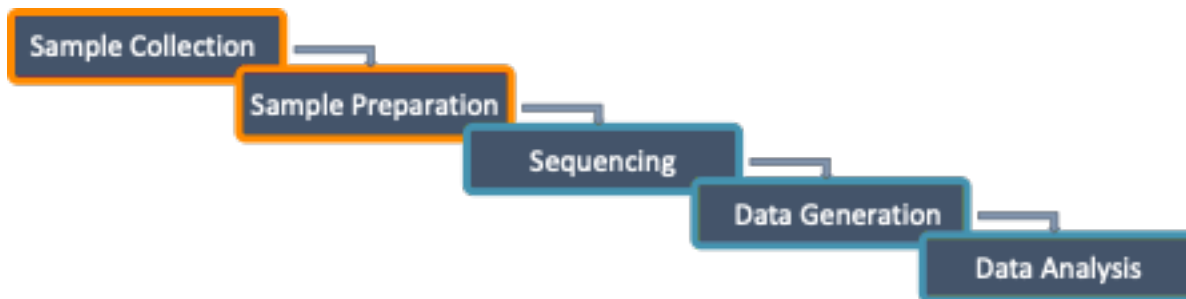
292 Some genes give rise to observable traits, termed phenotypes, like hair or eye color, blood type,
293 and facial features. Other phenotypes may include the presence of or susceptibility to certain
294 medical conditions, such as sickle cell anemia, cystic fibrosis, Huntington’s disease, forms of
295 muscular dystrophy, or certain cancers. Importantly, the genome reveals a great deal of
296 information about an organism as well as its relatives.

297 Genomic data is immutable, associative, and conveys important health, phenotype, and personal
298 information about individuals and their kin (past and future). In some cases, small fragments of
299 genomic data stripped of identifiers can be used to re-identify persons, though the vast majority
300 of the genome is shared among individuals [3][4][5][6][7][8].

301 2.1. Genomic Information Lifecycle

302 DNA sequencing, a key method researchers use to understand genomic information, identifies
303 the order of base pairs across a stretch of DNA or even the whole genome. This information
304 helps researchers understand genome organization and identify any similarities or differences
305 that may be present across organisms. DNA sequence data are generated by different
306 stakeholders, including academic researchers, healthcare providers, government agencies, and
307 industry. These entities may generate and store sequence data for internal use, make data
308 available for public access, distribute it to select entities or communities of interest, and/or
309 aggregate their results with others.

310 [Figure 1](#) summarizes the genomic lifecycle and contrasts those that are out of scope for this
311 report: Sample Collection and Sample Preparation with what is considered in scope: Sequencing,
312 Data Generation, and Data Analysis.



313 **Fig. 1.** Genomic Data Lifecycle; Naveed et.al. [6]

314 The genomic data lifecycle begins with sample collection, followed by DNA extraction and
315 several processing and quality control steps performed in the laboratory. Sample preparation
316 includes the processes where collected DNA is purified, fragmented, and bound to adapter
317 molecules to generate a library for sequencing. Libraries are collections of DNA molecules that
318 are compatible with the sequencing system. The prepared sample is then applied to a specially
319 designed substrate, like a glass flow cell, that is used by the sequencing system to read the order
320 of the DNA bases. Safeguarding sample collection and DNA processing requires applying
321 appropriate physical security as opposed to cybersecurity controls, and is deemed out of scope.
322 One exception is the collection of human DNA in the context of clinical care or study. In those
323 instances, human DNA sample collection and downstream steps are subject to additional privacy
324 considerations and informed consent.

325 Once a sample has been sequenced, the instrument outputs raw data. Resulting raw data file
326 formats differ across manufacturers of sequencing systems but typically contain identifiers, reads
327 (e.g., DNA sequences read by the instrument), and quality scores for each base call.

328 The resulting raw data are typically stored in a computer on-premises or in the cloud for
329 processing and analysis. Data processing and analyses can be performed by open-source or
330 commercially developed software or algorithms. The steps performed vary based on the intended
331 use of the sequence data. Example data processing steps include demultiplexing (if multiple
332 samples were sequenced within a single run), filtering out low-quality reads, aligning to a
333 reference genome, identifying variants, or annotating genomic variants (e.g., with how they may
334 change genes). The data may then be analyzed to identify, for example, variants across samples
335 or correlate variants with phenotypes. Common clinical applications of DNA sequencing include
336 diagnosing inherited diseases, targeting more effective drugs with precision medicine (e.g., for
337 cancer), and better understanding genetic risk factors for diseases.

338 These data may be made accessible to other researchers or uploaded to public or controlled
339 access databases, like the sequence read archive (SRA) or the database of Genotypes and
340 Phenotypes (dbGaP) that is run by the National Center for Biotechnology Information (NCBI) at
341 the NIH.

342 In healthcare, genomic data can be medically relevant at any time in a patient's life, thus
343 integration with an electronic health record (EHR) is essential. However, the nature and size of
344 the data presents a challenge with incorporating it into EHR workflows. A whole genome
345 sequence can be used to identify an individual in some cases, and the raw data can exceed 100
346 gigabytes in size. Additionally, as research advances, a genomic sequence may need to be re-
347 analyzed to provide the most up-to-date information for a patient's current medical condition.
348 Thus, the portability, privacy, chain-of-custody, interoperability, consent management, and re-
349 interpretation of genomic data are all key points for its effective use in healthcare.

350 Like other sensitive data, at each stage in the genomic data lifecycle from creation to storage and
351 analysis to dissemination, the data can be at risk of being intercepted, corrupted, overwritten, or
352 deleted.

353 **2.2. Next Generation Sequencing**

354 Since the introduction of automated sequencers in 1986 [9], sequencing costs have dropped by
355 several orders of magnitude [10]. Both the number and types of DNA sequencers have
356 proliferated for clinical and research applications. The first sequencing platforms used Sanger-
357 based chemistries, relying on the chain termination method, and were limited in the numbers of
358 sequences that could be read simultaneously. Advances in sequencing approaches and
359 chemistries have led to second- and third-generation sequencing platforms, often called next-
360 generation sequencing (NGS). NGS systems have been developed by several manufacturers with
361 varying read lengths and accuracy for different applications. Compared to automated Sanger
362 sequencing, which reads a single DNA molecule many times, NGS systems read many
363 molecules simultaneously, greatly increasing throughput and efficiency while reducing costs.

364 **2.3. Variant Calling**

365 Variant calling, the primary form of genomic analysis, is the identification of differences in the
366 genetic sequence between a sample and a reference genome. Variants can range from changes to
367 a single base at a given position in the genome to larger structural alterations, such as deletions,
368 insertions, duplications, inversions, and translocations. Variant identification is particularly
369 important in clinical medicine and molecular biology as alterations elucidate the role of genetic
370 elements in various diseases and facilitate our understanding of cellular processes.

371 Characterizing cancers is one area that employs variant identification. For instance, analyses of
372 somatic and germline genetic alterations—germline mutations are inherited by germ cells (sperm
373 and egg) during conception whereas somatic mutations occur after conception in non-germ
374 cells—yield insights to clinical outcomes and tumor origins. Somatic variants can be used to
375 target drugs that are more likely to be effective against the individual’s tumor.

376 As NGS technologies have progressed, so too have the computational pipelines for identifying
377 variants. Multiple software tools and workflows have been developed for different NGS
378 technologies and various variant calling applications. Precise variant calling requires having
379 complete and accurate reference genomes, as well as appropriate sequencing depth, accuracy,
380 read length, and analysis methods. Efforts including but not limited to the NIH Genome
381 Reference Consortium, Telomere-to-Telomere Consortium, and the NIST Genome in a Bottle
382 have significantly advanced the technical infrastructure, accuracy, and completeness of human
383 reference genomes and authoritatively characterize genomes to assess accuracy of variants
384 [11][12]. Recent systematic evaluations of germline and somatic variant calling pipelines
385 showcased the capabilities and limitations of several sequencing platforms and variant calling
386 software and workflows [13][14]. As sequencing systems continue to advance, variant calling
387 pipelines are likely to follow suit.

388 **2.4. Genome Editing**

389 New advances in DNA editing techniques, like clustered regularly interspaced short palindromic
390 repeats (CRISPR) and CRISPR-associated protein (CRISPR-Cas), promise precise, efficient, and
391 affordable ways to edit the genome through removing, adding, or altering DNA segments. As

392 next generation sequencing is helping scientists better understand the genomes of various
393 organisms, CRISPR-Cas systems are enabling researchers to modify genes more easily and
394 affordably than ever possible. When using genome editing in clinical applications to treat
395 diseases, NGS technologies are used to confirm the correct edit was made and identify any
396 unintended off-target edits. Although the net benefit of these technological advances is positive,
397 there exists an opportunity for misuse.

398 **2.5. Direct-to-Consumer Testing**

399 An alternative approach to measure DNA does not use sequencing but instead uses microarrays
400 of DNA probes to measure common variants at millions of specific genomic positions.
401 Microarrays are commonly used by DTC genetics companies for ancestry analysis and screening
402 for whether individuals may be carriers for particular disease-related variants. DTC genetic tests
403 are marketed to individuals without the involvement of a healthcare provider. These tests,
404 available from multiple companies, enable consumers to gain personalized insights into their
405 health and ancestry by examining their genetic profile. Results from consumers' genetic tests can
406 be accumulated in research databases to identify genetic sequence associations with certain
407 conditions. Moreover, third-party online genetic genealogy services allow consumers to upload
408 their genetic sequence information and identify related individuals.

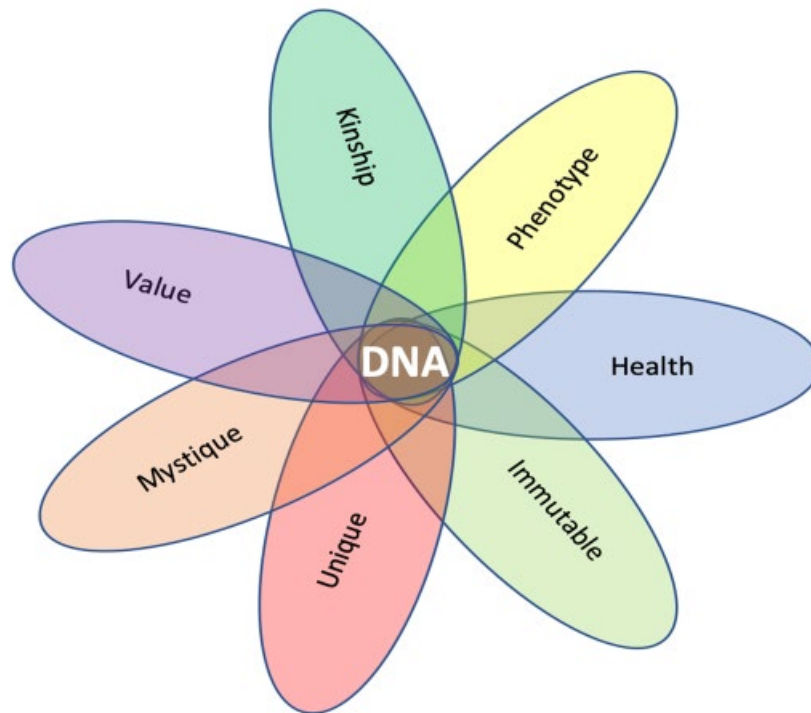
409 Advancements in microarray technology have allowed these companies to provide tests at
410 relatively affordable price points. The costs of DTC tests vary from company to company but
411 typically start at \$99. The DTC genetic testing market is estimated to be worth over \$1.3 billion
412 and is projected to grow to approximately \$3.5 billion by the end of 2026 [15].

413 **2.6. The Characteristics of Genomic Data**

414 Genomic data share attributes with other sensitive types of information, and as such, mirrors
415 their need for secure storage and transfer. Beyond these aspects there are several intrinsic
416 characteristics of genomic data. These concepts were discussed during the first NCCoE public
417 workshop held on May 18 and 19, 2021, and have been posited by some in the research
418 community [3]. [Figure 2](#) identifies seven features of genomic information that distinguish it from
419 other types of data. It is not any single characteristic, but instead the combination of these
420 intrinsic properties that highlight its value and sensitivity.

- 421 • **Phenotype.** Phenotype refers to the observable characteristics imparted by the genome,
422 such as size, appearance, blood type, and color. DNA can reveal a great deal of
423 information about an individual or their relatives, including phenotype or health
424 information.
- 425 • **Health.** Health means that DNA contains information about an organism's disease
426 presence, disease risk, vigor, and longevity. Clinical genetic testing can identify variants
427 within one's genome that may contribute to certain health outcomes.
- 428 • **Immutable.** Immutable means that an organism's DNA does not change significantly
429 during the organism's life. An individual's genome is practically immutable, with a

- 430 negligible lifetime mutation rate for most applications, which increases the long-term
431 consequences of a data breach.
- 432 • **Unique.** Unique means that individuals of species with sexual reproduction can be
433 identified. Except in the case of identical siblings, a person’s genome is unique to them.
 - 434 • **Mystique.** Mystique refers to the public perception about the mystery of DNA and its
435 possibly future uses.
 - 436 • **Value.** Value refers to the importance of the information content of DNA that does not
437 decline with time, but typically increases with time. Genomic information has value to
438 certain entities and that value is predicted to grow as we learn more about the genomes of
439 humans and other organisms.
 - 440 • **Kinship.** Kinship means that common ancestors and descendants of the organism can be
441 identified from DNA samples. Consumer genetic testing services provide information
442 about one’s ancestral lineage, including the potential to identify relatives.



443 **Fig. 2.** Characteristics of DNA (Adapted from Figure 1 of Naveed et al. [6])

444 **2.7. Balance Between Benefits and Risks for Uses of Genomic Information**

445 The U.S. research community, government, and private industry require genomic data sharing to
446 advance scientific and medical research and to maintain the country’s competitive advantage in
447 biotechnology. The transfer and sharing of genomic data are essential for understanding human
448 health, improving wellbeing, and accelerating scientific inquiry and advancements. For example,
449 in 2021 the NIH processed almost 40,000 requests for data access and has about three million
450 genotype microarray datasets and over 500,000 whole genome sequences [4][15]. Genomic data

451 enables precision medicine for rare diseases and cancer and is an enabler for CRISPR
452 technology. DTC genomic testing allows people to benefit from ancestry tracing, relative
453 matching, and health insights. In 2022, DTC companies reported more than 40 million
454 consumers have used their services [\[17\]](#). Genomic data can be useful for forensics to solve
455 crimes [\[18\]](#). Future uses of genomic data have further potential to advance the bioeconomy.

456 The genomic data transferred and shared represents tens of millions of individuals who provide
457 their information. In aggregate across all types of measurements, this data is touched by
458 thousands of entities (e.g., domestic, international, nonprofit, for-profit) that store, access,
459 manage, and use genomic and health-related data. These data sharing activities need adequate
460 technological and policy controls that allow research and enable commerce, as well as respect
461 the informed consent and privacy of the data subjects who expect protections from re-
462 identification [\[19\]](#).

463 Loss of control of genomic data can cause risks to privacy, personal security, and national
464 security, as adversaries can use genomic data for nefarious reasons such as surveillance,
465 oppression, and extortion. Genomic database breaches or other losses of data may result in thefts
466 of intellectual property and put the U.S. at a competitive disadvantage in biotechnology. As
467 reported by national security experts, security threats may arise through the creation of
468 population-specific bioweapons or compromised identities of national security agents [\[3\]\[20\]](#).
469 Cyber attacks have occurred on genomic databases, DNA sequencing instruments, and genomic
470 software tools [\[21\]\[22\]\[23\]\[24\]\[25\]](#).

471 3. Challenges and Concerns Associated with Handling Genomic Information

472 This section summarizes challenges and concerns for protecting genomic data in national
473 security, personal security and privacy, discrimination and reputational, economic, health
474 outcomes, and other potential future concerns. These challenges and concerns come from
475 discussions and presentations from the workshops and through literature review.

476 3.1. Potential National Security Concerns

477 As genomics research evolves, threats to national security continue to emerge. This section lists
478 potential national security concerns that were identified by workshop stakeholders. The
479 likelihood for these threats was not evaluated as part of the workshop.

480 The national security concerns identified include:

- 481 • **Infringement:** Genomic data can be used for population surveillance and oppression a
482 well as extortion of our citizens, military, and intelligence personnel [21].
- 483 • **Biological weapons:** Potential risks related to bioweapons directed against an
484 individual's DNA sequence or cloning an individual were deemed to be impractical and
485 not concerns that should drive genomics cybersecurity work [26].
- 486 • **Production of toxic products or infectious agents:** Peccoud J., et.al. [27], speculates
487 that due to the lack of integrity and tampering controls that typically exist, genomic data
488 could be corrupted by altering sequences or annotations and that, "These changes could
489 delay research programs or result in the uncontrolled production of toxic products or
490 infectious agents." This could result in significant loss of life as well as economic
491 consequences.

492 3.2. Privacy Challenges

493 Privacy challenges resulting from the use of human genomic data include problems for
494 individuals such as enabling blackmail, discrimination based on disease risk, and the revelation
495 of hidden consanguinity or phenotypes including health, emotional stability, mental capacity,
496 appearance, and physical abilities.

497 The U.S. intelligence community highlighted [3][21] concerns with the sparse regulation of the
498 genomic data collected and stored by DTC companies, along with the lack of recognition of U.S.
499 genomic data as an asset that needs export controls. Some researchers and consumer groups
500 suggested that due to the high reidentification risk, human genomic data should be consistently
501 classified as personally identifiable information (PII), while other researchers argued that this
502 would hinder research and the benefit to society [3]. It should be noted that regardless of the
503 classification, the privacy risks arising from processing genomic data remain.

504 Potential privacy problems identified include:

- 505 • **Re-identification of de-identified genomic data:** Human genomic data, even small
506 fragments of a person's whole genome, can usually be re-identified for some populations
507 when combined with available datasets, such as ancestry data, self-shared identified
508 genomic data of distant relatives, surname inference, age, etc. [6][7][8][28].

- 509 • **Unanticipated revelation of individuals' blood-relatives can lead to dignity loss**
510 **when those relationships are identified:** Consanguineal ties may be revealed that may
511 be embarrassing or incriminating, resulting in psychological or reputational harm.

512 3.3. Discrimination and Reputational Concerns

513 The discrimination and reputational concerns identified include:

- 514 • **False identification:** Sample mishandling, crime lab procedure deviations, or intentional
515 modification of the digital genomic data produce a risk to individuals being framed for
516 crimes they did not commit or extorted with the falsified genomic information.
- 517 • **Discrimination based on disease risk identified by an individual's genomic data:**
518 While some forms of discrimination in the U.S. are prohibited by law, the federal laws
519 are narrowly written and do not prohibit discrimination based on an individual's genomic
520 data for things such as life insurance, acceptance into the military, or senior residential
521 communities (e.g., based on Alzheimer's risk) [\[6\]\[29\]](#).
- 522 • **Unintended consequences from sample bias:** Artificial Intelligence (AI) and statistical
523 analysis techniques analyze large sets of genomic data to find diseases, risk factors for
524 diseases, make treatment decisions, and predict patient prognosis. Large datasets from
525 DTC and other samples of convenience in the U.S. are from predominantly those of
526 European descent. Minority communities are typically underrepresented. This bias of the
527 sample data can be amplified, particularly in AI techniques, and impact the results of
528 these analysis and prediction techniques that can result in discrimination [\[29\]\[30\]\[31\]](#),
529 resulting in potential harms to those whose genomic data are not represented in the
530 sample set. This bias of AI algorithms has been studied by NIST in the field of facial
531 recognition [\[32\]](#).

532 3.4. Economic Concerns

533 Economic concerns identified include:

- 534 • **Intellectual property infringement:** Exfiltration of genomic data can result in loss of
535 intellectual property for institutions, resulting in economic losses in the future.
- 536 • **Operational disruption:** Disruptions to the generation, storage, or use of genomic data
537 can negatively affect production operations, impacting agricultural food production,
538 biopharmaceutical output, and biomanufacturing.
- 539 • **Extortion:** Bad actors or nation states may extort individuals based on their genomic
540 data, threatening to reveal sensitive health or kin information encoded in the individual's
541 genome [\[21\]](#), resulting in financial, psychological harm, and reputational loss to the
542 individual.
- 543 • **Penalties and liabilities:** Negligence in the cybersecurity controls of genomic data by the
544 company or institution could violate their regulatory obligations, jeopardize their access
545 to data in the future, and result in significant financial loss through imposed penalties.
- 546 • **Revoke data access:** Negligence in adequate security or privacy protection can cause
547 individuals to no longer consent to have their genomic data used in scientific studies,
548 which could reduce the effectiveness of research and threaten the bioeconomy.

549 **3.5. Health Outcome Concerns**

550 A patient’s genomic data needs to be effectively incorporated into their EHR to allow for
551 effective clinical care. Concerns in healthcare include portability, chain-of-custody, re-
552 interpretation of genomic data, and consent management.

553 **3.6. Other Potential Future Concerns**

554 Other potential future concerns identified include:

- 555 • **Loss of individuality:** While not a realistic threat at present or in the near future, cloning
556 of an individual could theoretically be accomplished with their genomic information.
557 This could result in psychological or financial harm to the person and their clone [33].
- 558 • **Human engineering:** Genomic information could be used for eugenics, more finely
559 targeting kinship or other human traits [34]. This could result in societal harm or
560 psychological harm to the individual.

561 **3.7. Summary of Challenges and Concerns with Genomic Data**

562 The U.S. research community, government, healthcare, and private industries handle genomic
563 data and require genomic data sharing to advance scientific and medical research, improve health
564 outcomes, and maintain the country’s competitive advantage in biotechnology. Genomic data
565 often needs to be aggregated from multiple studies to address pressing research questions, but
566 challenges such as the differing subject consents used in the different studies can present
567 difficulties that need adequate technological and policy controls that respect the informed
568 consent and privacy of the data subjects [3]. Though it can be time-intensive and difficult,
569 responsible data sharing and analytics facilitates commerce, research, and healthcare outcomes
570 while protecting subject privacy against re-identification and respecting subject informed
571 consent. Further, additional guidance is needed to promote the safe and secure sharing of
572 genomic data while addressing threats to national security, personal security and privacy,
573 reputations and civil liberties, the economy, health, and the future of humanity.

574 **4. Current State of Practices**

575 This section identifies many of the cybersecurity, privacy, and risk management practices issued
576 by U.S. Government, industry, and international entities. While cybersecurity and privacy
577 practices are inter-related and support overall risk management practices, this section identifies
578 practices specific to each area. Subsections highlight the NIST Risk Management Framework
579 (RMF), the NIST Cybersecurity Framework, and the NIST Privacy Framework, along with
580 legislation, frameworks, alliances, and other resources that can be leveraged to improve the
581 protection of genomic data. The final subsection identifies guidance, technical, and
582 policy/regulatory gaps that can then be addressed by solutions proposed in the subsequent
583 section.

584 **4.1. Risk Management Practices**

585 Risk management is the process of managing risks to organizational assets and operations,
586 individuals, and other entities (e.g., nations). The NIST RMF [35] provides a well-established
587 method for information security risk management that can apply to any system or organization.
588 Federal risk management related initiatives such as Executive Order (EO) 14028 [36] and the
589 subsequent guidance published by NIST define methods for protecting the software supply
590 chain. The “Framework for Responsible Sharing of Genomic and Health-Related Data” [37]
591 provides an additional resource from the international community.

592 **4.1.1. U.S. Government Resources**

593 For the U.S. Federal government, the Federal Information Security Management Act (FISMA
594 2002) served as the initial driver for cybersecurity risk management programs. The Office of
595 Management and Budget (OMB) Circular A-130, “Managing Federal Information as a Strategic
596 Resource” requires Executive agencies to leverage NIST guidance. The NIST RMF and the
597 Federal Risk and Authorization Management Program (FedRAMP) are examples of risk
598 management processes used by federal agencies.

599 **FISMA** required each federal agency to develop, document, and implement an agency-wide
600 program to provide information security for the information and systems that support the
601 operations and assets of the agency, including those provided or managed by another agency,
602 contractor, or other sources. The Federal Information Security Modernization Act (FISMA 2014)
603 [38] includes updates to address evolving cybersecurity concerns, reduce reporting burdens,
604 strengthen continuous monitoring in systems, and reporting of incidents.

605 **OMB Circular A-130** requires executive agencies within the federal government to plan for
606 security, ensure that appropriate officials are assigned security responsibility, periodically review
607 the security safeguards in their systems, and authorize system processing prior to operations and
608 periodically, based on risk.

609 **The NIST RMF** (defined in NIST SP 800-37 Revision 2) [35] provides a structured, yet
610 flexible, process for managing cybersecurity and privacy risk that includes steps for preparation,
611 system categorization, control selection, control implementation, control assessment, system
612 authorization, and continuous monitoring. Risk management involves more than complying with

613 regulations or technical controls and should be tailored to each organization’s mission,
 614 regulatory environment, and risk tolerance. [Table 1](#) NIST Risk Management Framework
 615 Discussion for Genomic Data provides an overview of the RMF and a brief description of its
 616 relevance to genomic data.

617 **Table 1.** NIST Risk Management Framework Discussion for Genomic Data

RMF Step	Overview	Genomic Data Relevance
Prepare	The organization carries out essential activities to prepare for risk management. This includes identifying key risk management roles, establishing an organization-wide risk management strategy and risk tolerance, assessing organization-wide security and privacy risks, and developing an organization-wide continuous monitoring strategy. At the system level, the organization must understand information types in the system, conduct a system-level risk assessment, identify the relevant requirements—including the applicable federal and state regulatory requirements.	It is important that all the relevant regulatory requirements are considered at a state or national level depending on the subject’s citizenship, where the data was collected, and where it is being processed. An important feature of human genomic and health data that organizations must prepare for is that the informed consent may have specific restrictions on how that data is used and that the informed consent often differs depending on when and under what circumstances the data was collected.
Categorize	Using outputs of the Prepare step, categorize the system and information processed, stored, and transmitted and gauge the potential loss of the Confidentiality, Integrity, and Availability (CIA) triad for the information, system, and processes.	Data categorization helps an organization identify the appropriate controls to properly mitigate the risk to an acceptable level. Human health and genomic data may additionally be sub-categorized by the allowed uses of the subject’s informed consent which may restrict how it is processed.
Select	Controls are selected to manage risk in accordance with the risk management strategy and within the risk tolerance levels.	Specific controls related to genomic data have yet to be identified but could be used as guidance for organizations managing related risks.
Implement, Assess, and Authorize	The system security plans are updated to reflect the implemented state of the controls. Controls are assessed (which includes validation and verification) to ensure they are in place, operating as intended, and achieving the desired results. The Authorize step requires a senior official to understand the residual risks and agree that they are acceptable to the organization before the system is put into operation (or continues to operate).	These steps enable the organization to quantify and characterize the risks associated with the current protections implemented (or not implemented) for the genomic data. Organizations would then manage inherent risk or residual risk through plans of action and milestones (POA&Ms) identifying the resources and timelines required to address residual risks.
Continuous Monitoring	Once the system is operational, an organization must ensure the system is operating as intended and within the acceptable risk tolerance of the organization. Subject matter experts with appropriate expertise are necessary at this stage. Periodic evaluation is needed to keep up with changes as they happen, including new threats, regulatory changes, and technology changes. Organizations must also consider end-of-life procedures for their systems and data.	Organizations must monitor both the relevant threats and outstanding vulnerabilities to determine the risk posed to the organization. Response and recovery plans must be in place to appropriately address events and incidents as they occur to minimize exposure, protect the genomic data, inform users of breaches, and recover from incidents.

618 **FedRAMP** [39] is a risk-based approach, aligned to the NIST RMF, for protecting cloud-based
619 services and systems and includes continuous monitoring and an independent assessment
620 requirement. FedRAMP is governed by the OMB, the General Services Administration
621 FedRAMP Program Management Office (PMO), and a FedRAMP Joint Authorization Board.
622 FedRAMP outlines guidance for securing a cloud-hosted system and FedRAMP authorization
623 may be re-used across multiple government agencies. Federal agencies can use a FedRAMP
624 system as part of their overall solution for implementing FISMA requirements.

625 **4.1.2. U.S. Government Initiatives**

626 President Biden’s administration issued EO 14028, Improving the Nation’s Cybersecurity [36]
627 on May 12, 2021, to direct federal agencies to enhance cybersecurity risk management practices.
628 In response to the EO, NIST issued guidance on critical software and updated guidance on
629 protecting the software supply chain [40]. Genomic sequencing environments can leverage this
630 guidance to implement supply chain protections that will improve visibility into provenance and
631 related software components. More recently, EO 14081, Executive Order on Advancing
632 Biotechnology and Biomanufacturing Innovation for a Sustainable, Safe, and Secure American
633 Bioeconomy [1] requires identifying cybersecurity risks in biotech.

634 NIH continues to play a leading role in genomic research and the protection of genomic data.
635 Examples of NIH-led efforts include:

- 636 • NIH manages dbGaP, with almost 40,000 requests for access to datasets within this
637 database (2022) that must be shared responsibly.
- 638 • NIH is leading genomic-related efforts to move from identity-based access to authority-
639 based access. NIH recommends moving to authority-based access, which focuses on what
640 a researcher and the software are authorized to do within a system, rather than the identity
641 of the researcher. In continually monitoring changing threats, regulations, and
642 technology, NIH has determined that identity-based security models have become
643 insufficient to adequately mitigate the risk for their mission. Identity-based security
644 models lack context and have poor federation. Once a user has been authenticated, there
645 are no further limitations on what the user (or the software used by the user) is authorized
646 to do with the data. Authority-based access uses tokens and is more consistent with the
647 Zero Trust Architecture (ZTA) model and the principle of least privilege [41][42].
- 648 • The NIH has been working with the Global Alliance for Genomic Health (GA4GH) to
649 ensure that their implementation of authority-based access is compatible with the
650 GA4GH Passport for researcher authorization.
- 651 • NIH has a significant concern about the confinement problem for genomic data sharing—
652 preventing an authorized user from sharing data with others—and continues to look for
653 solutions to this issue. The NIH addresses the confinement problem with contractual
654 controls, as the currently available technical controls have limitations that are not
655 compatible with the NIH mission. However, technical controls that address the
656 confinement problem and have limitations that are more compatible with the NIH
657 mission are being researched.

658 **4.1.3. International Resources and Regulations**

659 The “Framework for Responsible Sharing of Genomic and Health-Related Data” [37] is an
660 international framework specific to genomic data risk management. It addresses international
661 data sharing, collaboration and good governance concerns, and the protection and promotion of
662 the welfare, rights, and interests of individuals from around the world in genomic and health-
663 related data sharing.

664 In the European Economic Area (EEA), organizations must consider the “Schrems II” ruling
665 [43][44] which addresses confidentiality and data residency of data provided by citizens in the
666 member states of the EEA. In Asia, India is considering expanding its data protection law in a
667 manner similar to the European Union General Data Protection Regulation [45][46], and China is
668 compelling Chinese companies to share data they have collected with the government [21].
669 Further, China and Russia severely restrict the sharing of genomic information [47][48][49][50].

670 **4.2. Cybersecurity Best Practices**

671 Cybersecurity practices support overall risk management. This section identifies federal,
672 international, and industry resources that organizations can use to help identify, prioritize, and
673 implement cybersecurity capabilities. The NIST Cybersecurity Framework outlines categories of
674 capabilities and provides an overall methodology for prioritizing cybersecurity investments.
675 Another widely used federal resource is the guidance on implementing a ZTA. Additional
676 practices introduced in this section include information sharing and analysis centers, software
677 supply chain management, and the GA4GH Data Security Infrastructure Policy [51].

678 **4.2.1. U.S. Government Resources**

679 The U.S. Government publishes multiple resources that support cybersecurity practices. This
680 section describes 1.) the NIST Cybersecurity Framework, 2.) ZTA guidance, and 3.) Benchmarks
681 or Security Technical Implementation Guides (STIGs), produced by the Defense Information
682 Systems Agency (DISA).

683 **The NIST Cybersecurity Framework** [52] is a voluntary, comprehensive framework
684 developed by the federal government, that is applicable to any organization (federal, academic,
685 private sector, etc.) wanting to improve their cybersecurity risk management. It is very broad, not
686 genomic specific, and can be tailored for real-world deployments within specific industries.

687 **ZTA** is an evolving cybersecurity paradigm described in NIST SP 800-207 [53] that identifies
688 target requirements and capabilities that align with genomic data protection goals. ZTA enables
689 secure authorized access to resources individually in situations when many users need access
690 from anywhere, at any time, from any device to support the organization’s mission. Data is
691 programmatically stored, transmitted, and processed across different boundaries under the
692 control of different organizations to meet ever-evolving business use cases. In a ZTA, no implicit
693 trust is granted to assets or user accounts solely on their physical or network location. Allowing
694 assets and users to only access the required resources for fulfillment of their role in the
695 organization’s mission provides for defense-in-depth.

696 **DISA STIGs** [\[54\]](#) and similar benchmarks can define secure configurations for cyber-physical
697 systems with operating systems (such as sequencers) that can be assessed and maintained. These
698 guides provide hardening and defense in depth, greatly improving an organization’s security
699 posture as the default settings of operating systems are typically optimized for usability and have
700 many insecure settings not appropriate for securing high-value data. Security benchmarks
701 support the secure configuration of hardware and software throughout the genomic lifecycle.

702 **4.2.2. International Resources**

703 The **Bioeconomy Information Sharing and Analysis Center (BIO-ISAC)** [\[55\]](#), an
704 international non-profit organization, addresses threats unique to the bioeconomy, sharing threat
705 intelligence among the stakeholder community. BIO-ISAC facilitates the responsible disclosure
706 of detected or suspected vulnerabilities found in sequencers. This collaborative approach to
707 cybersecurity complements individual organizational efforts. BIO-ISAC is an example of an
708 industry stakeholder providing cybersecurity services, such as emergency threat hunting for
709 organizations who are under active attack, to bioeconomy organizations who choose to augment
710 their own resources.

711 The GA4GH has developed a **Data Security Infrastructure Policy** [\[56\]\[57\]](#) that is a “set of
712 recommendations and best practices to enable a secure data sharing and processing ecosystem.”
713 They provide frameworks and standards for responsible genomic data sharing. While the
714 proposed frameworks and standards provide foundational principles for genomic data sharing,
715 they are high-level guidelines and lack articulation of in-depth and comprehensive security
716 features and solutions that are necessary for a real-world development and deployment of such
717 principles.

718 **4.2.3. Industry Resources**

719 Industry cybersecurity practices including software bill of materials (SBOM), automated
720 vulnerability scanners, and the use of manual expert analysis support the protection of genomic
721 data. These industry practices are becoming more widely implemented even in the U.S. federal
722 government and should be considered for all systems processing genomic data.

723 **SBOMs** are widely used across industry and now being promoted by the Cybersecurity and
724 Infrastructure Security Agency (CISA) [\[58\]](#) as an essential building block in software security
725 and software supply chain risk management. SBOMs contain a nested inventory identifying the
726 ingredients that make up software components and provide transparency to cybersecurity
727 professionals and users to facilitate effective patch management and mitigation of newly
728 discovered software vulnerabilities.

729 **Vulnerability scanners** are automated tools that examine software systems for
730 misconfigurations and software flaws that compromise the cybersecurity of the system. They are
731 highly useful in detecting unintended mistakes in threat mitigation and are an important part of
732 assessing a system for cybersecurity. Additionally, ZTA guidance from OMB M-22-09 requires
733 federal agencies to conduct manual expert analysis of systems by approved external third-party
734 testing capabilities.

735 **4.3. Privacy Best Practices**

736 This section describes the current state of privacy practices related to genomic data including
737 U.S. Government and industry resources along with relevant regulations.

738 **4.3.1. U.S. Government Resources**

739 U.S. Government privacy resources include the NIST Privacy Framework, Federal laws
740 including Genetic Information Nondiscrimination Act of 2008 (GINA) [59], Health Insurance
741 Portability and Accountability Act (HIPAA) [60] and the Common Rule.

742 **The NIST Privacy Framework** [2] is a voluntary framework that helps organizations identify
743 and manage privacy risk within an organization’s broader enterprise risk portfolio. Like the
744 NIST Cybersecurity Framework, the NIST Privacy Framework can be used across federal and
745 non-government organizations. While not specific to any use case, it can be tailored to address
746 genomic data privacy requirements for those producing, storing, or processing genomic data.

747 **GINA** [59] includes two key provisions that prohibit group insurance providers and employers
748 of more than 15 people from using genetic information to discriminate against individuals.

749 **HIPAA** [60] protects patient confidentiality by placing restrictions on sharing protected health
750 information (PHI) by HIPAA-covered entities (i.e., insurance companies and healthcare
751 providers) [61]. A 2013 amendment modified the HIPAA Privacy Rule to consider genetic
752 information as PHI [61]. Non-covered entities, like employers, law enforcement, life insurance
753 companies, school districts, and state agencies, are not required to comply with HIPAA
754 protections. Importantly, the HIPAA Privacy Rule does not place restrictions on using or
755 disclosing PHI of de-identified data. For defining whether data is sufficiently de-identified for
756 disclosure, the U.S. Department of Health and Human Services (HHS) has defined two
757 acceptable methods [61]. One is expert determination § 164.514 (b)(1), which requires applying
758 statistical or scientific principles and is used by the Department of Veterans Affairs (VA) [3]
759 when handling genomic data. The second is the “Safe Harbor Method” § 164.514 (b)(2), in
760 which data is considered deidentified if 18 types of identifiers are removed and there is “no
761 actual knowledge that the residual information can identify an individual.” However, the second
762 method is inappropriate for genomic data as genomic data is of extremely high dimensionality
763 and even small fragments can re-identify individuals with a high degree of probability with
764 current technology and publicly available information [3].

765 **The Common Rule** [62], officially known as the Federal Policy for the Protection of Human
766 Subjects, establishes a standard of ethics for government-funded human subjects research. A
767 2017 update required all federally funded human subjects research to obtain meaningful
768 informed consent. Study participants must be informed of how their genomic information will be
769 used, who may access it, and what are the potential risks associated with release of PHI [3].
770 Human subjects participating in clinical studies funded by the NIH are automatically granted a
771 Certificate of Confidentiality that safeguards their PHI. The Certificate of Confidentiality
772 compels investigators and institutions to withhold PHI from civil or criminal proceedings. These
773 certificates aim to promote clinical study participation by assuring a certain level of privacy.
774 Non-federally funded studies, however, are not obligated to provide Certificates of
775 Confidentiality.

776 **4.3.2. Industry Resources**

777 **Future of Privacy Forum** (FPF) Best Practices [56] is an industry supported voluntary set of
778 principles targeted at the DTC genomic testing market and uses a Fair Information Practice
779 Principles (FIPPs) based framework. It explicitly recognizes that genomic data absent other
780 identifiers can usually be re-identified and thus continues to need strong protection provided by
781 technical or contractual controls. The framework lists several important privacy issues and uses
782 the FIPPs principles to address them at a high-level. It lacks specifics on cybersecurity and risk
783 management of genomic data and is limited in focus to DTC genetic testing companies.

784 **4.3.3. International Resources**

785 **Global Alliance for Genomics and Health: Data Privacy and Security Policy** [57] is a
786 document focused on the responsible data sharing of genomic and health-related data. It sets
787 forth some general policies based on the four principles of (1) Respect Individuals, Families, and
788 Communities, (2) Advance Research and Scientific Knowledge, (3) Promote Health, Wellbeing,
789 and the Fair Distribution of Benefits and (4) Foster Trust, Integrity, and Reciprocity. It is an
790 international guideline and does not explicitly consider U.S. privacy regulations.

791 **4.4. Summary of Gaps in the Protection of Genomic Data**

792 The NCCoE engaged with the public to identify common challenges, what is unique about
793 genomic data, solutions that are available, and gaps faced by those who produce, store and
794 process genomic data. Subject matter experts with experience in bioeconomy, cybersecurity,
795 privacy, sequencing technologies, cloud hosting, and computation of genomic data were
796 contacted. A five-and-a-half-hour virtual workshop was held on January 26, 2022, with nearly
797 500 stakeholders to understand the threats related to the privacy and security of genomic data
798 and identify gaps in protection. Findings from that workshop were presented at the Healthcare
799 Information and Management Systems Society (HIMSS) Global Health Conference & Exhibition
800 on March 16, 2022, where feedback from attendees was solicited. An additional, second virtual
801 workshop was held on May 18 and 19, 2022, where possible solutions were explored, and
802 additional gaps were identified.

803 This section summarizes the gaps identified in guidance, technical solutions, and
804 policy/regulations.

805 **4.4.1. Guidance Gaps**

806 Current cybersecurity and privacy risk management guidance does not address the specific and
807 unique requirements of genomic data. This highlights the opportunity to address the following
808 gaps:

- 809 • Publish voluntary guidance for protecting genomic data through the development of
810 NIST Cybersecurity Framework and NIST Privacy Framework Profiles.

- 811 • Provide expertise to help life-sciences researchers leverage NIST guidance that supports
812 FISMA and related requirements to improve the cybersecurity and privacy risk
813 management of genetic data.
- 814 • Publish cybersecurity benchmarks or STIGs for the secure configuration of commercially
815 available sequencers and the associated data analysis pipelines to address the current gap
816 in visibility into how sequencer operating system and software settings have been
817 configured for system hardening and the ability to assess these configurations after
818 maintenance or updates.

819 **4.4.2. Technical Solution Gaps**

- 820 • Currently, sequencer manufacturers do not provide SBOMs for their devices. Therefore,
821 security professionals have no visibility into potential vulnerabilities of their software and
822 thus cannot adequately advise users of sequencers on how to address discovered software
823 vulnerabilities through patching or other mitigation measures.
- 824 • Sequencers are typically connected to a network and the internet. This provides access to
825 the manufacturer for updates and transfer of files to secure storage. There is no guidance
826 on the network addresses and the corresponding network protocols that are required for
827 sequencers to effectively operate. If this guidance existed, it would be possible to
828 microsegment sequencers on the network, providing them with only the required
829 resources for their proper functioning in keeping with ZTA cybersecurity principles. This
830 would mitigate the possibilities of adversaries exploiting vulnerabilities in the
831 sequencer's hardware or software for exfiltration of data as well as using the sequencers
832 as entry points that allow lateral movement throughout the enterprise network.
- 833 • The general problem of data confinement, (i.e., authorized users and/or their software
834 sharing unauthorized access to data), is a well-known unsolved problem in cybersecurity.
835 Due to the privacy risks to subjects as well as the high value of many types of genomic
836 data, the confinement problem is of relevance to genomic data. This problem is most
837 commonly addressed by contractual controls that can be particularly complex when the
838 controls are between multiple organizations. Further, contractual controls typically do not
839 prevent unauthorized data sharing, but provide penalties if it is done. These penalties
840 typically cannot redress the privacy loss of patients and data subjects.
- 841 • Most genetic data sharing and processing occurs in cloud environments, frequently
842 leveraging containers (e.g., Docker or Pods). Many cybersecurity vulnerability scanners
843 are not optimized for scanning containers, resulting in an inability to identify certain
844 vulnerabilities and a high number of false positives.
- 845 • In the healthcare of a patient, the nature and size of genomic data presents a challenge
846 with incorporating it into EHR workflows. Fast Healthcare Interoperability Resources
847 (FHIR) as a genomic cybersecurity solution provides important security features, such as
848 tracking provenance and consent as well as providing encryption for privacy. However,
849 additional work is needed to scale for genomic data size files.

850 **4.4.3. Policy/Regulatory Landscape**

- 851 • While some forms of information are subject to export controls, U.S. laws do not treat
852 genetic data as a national security asset, but primarily focus on privacy and intellectual
853 property (IP) protection. Few restrictions prevent a U.S. company from selling genetic
854 data to parties outside the U.S. Some of these gaps may be addressed through efforts to
855 implement EO 14081.
- 856 • Data held by DTC genetic testing companies are not subject to HIPAA or privacy and
857 security requirements that apply to health care providers, as consumers send samples
858 directly to the companies without the involvement of a healthcare provider. Only recently
859 have several states enacted consumer protection laws directed at DTC genetic testing
860 companies. These laws provide consumers with additional rights such as requiring
861 consent for data sharing or granting consumers the ability to access or delete their data
862 [\[63\]](#)[\[64\]](#)[\[65\]](#).
- 863 • Because laws of some countries may prevent the aggregation of a global human genome
864 representing all people groups [\[45\]](#)[\[48\]](#), AI methods trained on genomic datasets may be
865 biased with respect to U.S. citizens whose heritage includes portions of the under-
866 represented people groups.
- 867 • Multiple peer-reviewed studies [\[6\]](#)[\[7\]](#)[\[8\]](#)[\[28\]](#) have demonstrated that even small parts of
868 genomic data can be re-identified with high probability of success. The NIH and the VA
869 consider the re-identification risk of genomic data and have department-specific policies
870 to prevent this occurrence [\[3\]](#). Existing de-identification approaches, like the HIPAA safe
871 harbor provision, do not fully circumvent re-identification risk because genomic data
872 alone may be sufficient to re-identify an individual. Moreover, HIPAA only applies to
873 covered entities and genomic information in the context of PHI.

874 **5. Available Solutions to Address Current Needs**

875 This section describes opportunities to address the gaps identified in Section 4 through
876 technologies, processes, and guidance.

877 **5.1. NIST Cybersecurity Framework Profile for Processing of Genomic Data**

878 Genomic data is highly valuable to organizations in the bioeconomy. Organizations need to
879 apply appropriate cybersecurity capabilities for the protection of the data. The NIST
880 Cybersecurity Framework is a voluntary, comprehensive, and applicable framework that
881 organizations can use as part of an overall risk management plan. However, the NIST
882 Cybersecurity Framework is a general framework, and many industries have found it beneficial
883 to have a Cybersecurity Framework Profile based on that industry's mission objectives.

884 An industry-specific Cybersecurity Framework Profile provides several benefits. A Profile can
885 be used as a standardized approach for preparing a cybersecurity plan appropriate for the security
886 of genomic data; a method to select controls and mitigations appropriate for the high value of
887 genomic data; and can be used for gap analysis for organizations already storing and processing
888 genomic data to examine their cybersecurity posture against and prioritize upgrades to that
889 posture. Thus, a Cybersecurity Framework Profile specifically tailored to the mission objectives
890 of those who handle genomic data would be a valuable addition for securing the bioeconomy.
891 The NCCoE engaged stakeholders from the genomic community to create the Cybersecurity
892 Framework Profile for genomic data and anticipates publishing the Profile in the summer 2023.

893 **5.1.1. Use Case Description**

894 Organizations that handle genomic data need to protect that data due to both its high value as
895 well as the privacy risk to individuals if the data were exposed. The organization needs to select
896 appropriate controls to reduce the risk to the confidentiality, integrity, and availability of the data
897 to an acceptable level. An organization must develop an appropriate cybersecurity plan to
898 accomplish the task in a cost-effective manner so that the mission objectives of the organization
899 can be achieved. After implementing the plan, the organization needs to periodically assess its
900 cybersecurity posture considering new technology and threats and to analyze their current
901 posture versus an appropriate target posture to determine if there are gaps in its cybersecurity
902 posture and prioritize the remediation of those gaps.

903 **5.1.2. Solution Idea**

904 There exist a number of Target Profiles for specific use cases for the NIST Cybersecurity
905 Framework [66], however none address the bioeconomy, considering its challenges and mission
906 objectives. A NIST Cybersecurity Framework Profile for genomic data would be particularly
907 relevant because of the high value of genomic data and the important risks to the nation,
908 economy and individuals' loss of that data can have. The Profile could highlight controls that
909 address the unique aspects of genomic data and the most important considerations for
910 organizations that process genomic data.

911 **5.1.3. Expected Benefits**

912 A Cybersecurity Framework Profile that is specific to genomic data would assist organizations
913 who aggregate and process genomic data by highlighting gaps in their cybersecurity capabilities.
914 It would also assist organizations in prioritizing and implementing additional capabilities or
915 controls.

916 **5.2. NIST Privacy Framework Profile for Processing Genomic Data**

917 Human subjects who provide their genomic data for research expect that their privacy will be
918 protected by all organizations involved in the research process. Human genomic data can reveal
919 a great deal of personal information, such as physical traits, predisposition to certain health
920 conditions, and biological relationships. Individual privacy may be impacted when an individual
921 is identified, and other protected or sensitive information is uncovered or made available. For
922 instance, identifying an individual in a genomics database for patients with a particular medical
923 condition may impact the privacy of that individual by revealing sensitive health information.
924 The genetic data itself can serve as an identifier and provide additional information about the
925 contributor. Lastly, aggregate data in large genomic databases can lead to the identification of
926 individuals or their biological relatives [8][68]. Deidentifying genomic data is impossible
927 without destroying some or all of utility of that data [3].

928 In addition to a broad range of privacy considerations, human subjects provide their data under
929 informed consent, which may further limit the processing of their data to specific uses. Genomic
930 data collected with differing informed consents is often aggregated, but when this is done, care
931 must be taken that all data is still used within the boundaries specified in the informed consent of
932 each data subject. Organizations that process genomic data need to be able to effectively
933 communicate internally and externally about managing these and other privacy risks.

934 The NIST Privacy Framework [52], a voluntary tool, helps organizations to prioritize the policies
935 and technical capabilities they need to manage the privacy risks that may arise from data
936 processing, including processing genomic data. Like the Cybersecurity Framework, the Privacy
937 Framework can be applied to and tailored for the mission objectives for processing genomic data
938 to manage risks to individuals as well as related risk that can arise to organizations when
939 developing their products, systems, and services (e.g., reputational risk, loss of trust, financial
940 risk).

941 **5.2.1. Use Case Description**

942 Organizations that process human genomic data need to protect individuals from privacy risk.
943 These organizations also need to process their genomic data within the boundaries allowed by
944 each subject's informed consent and usually must respect and manage multiple different
945 informed consents. They must have a plan for effective technical controls and processes that will
946 reduce the privacy risk to an acceptable level while still accomplishing the mission objective in a
947 timely and cost-effective manner.

948 **5.2.2. Solution Idea**

949 There exist a number of Target Profiles for specific use cases for the NIST Cybersecurity
950 Framework [66]; however, no NIST Privacy Framework Profiles have been developed. A NIST
951 Privacy Framework Profile would be particularly relevant for genomic data because of the high
952 sensitivity of genomic data and the unique reidentification risks associated with genomic data.
953 The Profile could address the unique aspects of genomic data and highlight the most important
954 considerations for aggregators and processors of genomic data.

955 **5.2.3. Expected Benefits**

956 A Privacy Framework Profile that is specific to genomic data would assist organizations who
957 aggregate and process genomic data by highlighting gaps in capabilities intended to protect the
958 privacy of individuals. It would also assist organizations in prioritizing implementation of
959 additional privacy capabilities or controls.

960 **5.3. Automatic Network Micro-Segmentation of Sequencers with MUD**

961 The NCCoE recently developed model implementations of a solution to help secure internet of
962 things (IoT) devices that may have applications for genomic sequencers. NIST released NIST SP
963 1800-15, Securing Small-Business and Home Internet of Things (IoT) Devices: Mitigating
964 Network-Based Attacks Using Manufacturer Usage Description (MUD) in 2021 [68]. MUD
965 restricts a managed device to communicate with only an “allowlist” of permitted network
966 addresses (internet protocols and ports) specific to each type of device, consistent with ZTA
967 cybersecurity principles.

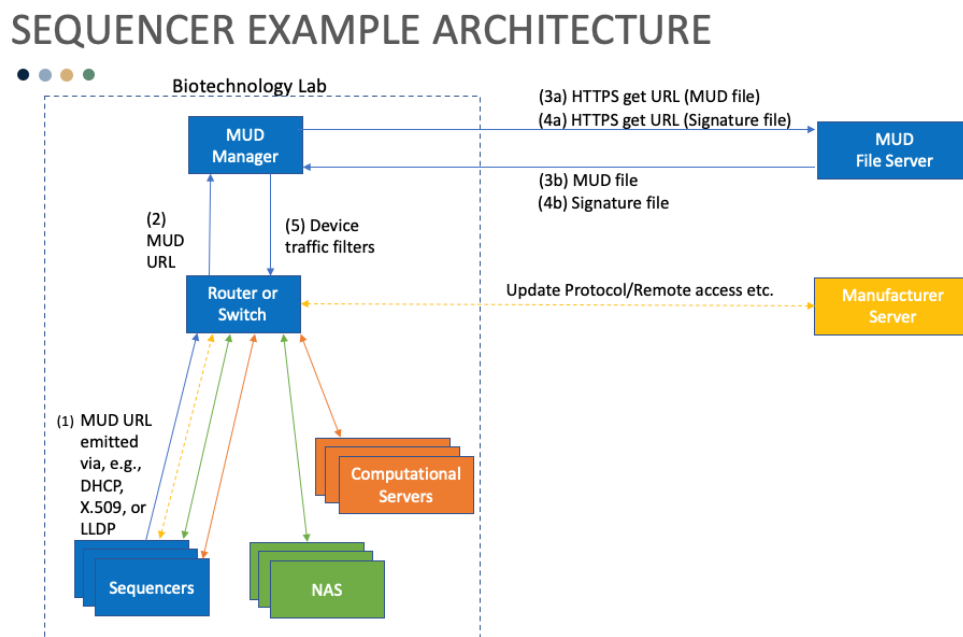
968 **5.3.1. Use Case Description**

969 Sequencers are expensive devices that are operated by custom software that provides only
970 limited visibility into network connections and configurations. Sequencers must be connected to
971 the internet and intra-network for manufacturer service, user access, and data storage, but the
972 number of communication addresses and protocols they require for proper operation are limited.
973 In these ways, sequencers can be compared to IoT devices [3]. Because of sequencers’
974 connectivity, they may be an attractive target for data exfiltration, ransomware deployment,
975 malware implantation, and lateral movement within the user’s enterprise network.

976 **5.3.2. Solution Idea**

977 MUD could be applied to sequencers. Since sequencers have a small number of communication
978 patterns, MUD enables a network to limit sequencer communication within the local network
979 and externally only to those resources needed for proper functionality. There is a MUD
980 architecture that implements MUD and contains a MUD Manager that provides micro-
981 segmentation of the device based on the manufacturer-provided allowlist. This micro-
982 segmentation provided by the MUD Manager greatly inhibits adversaries’ ability to access a
983 managed device; and if they do access a device, their ability to move laterally across the rest of

984 the network is severely restricted by the MUD Manager. [Figure 3](#) illustrates a potential MUD
985 architecture for a sequencer [\[4\]](#).



986 **Fig. 3.** Notional MUD architecture for a sequencer

987 MUD can be applied in partnership with the manufacturer. Existing tools can help either a
988 manufacturer or a third party develop an allowlist implemented using MUD for a device that
989 controls a device’s network traffic. A significant advantage of the MUD solution is that it can
990 also be implemented with legacy devices (whose software cannot be modified) to provide the
991 needed MUD information when connected to the network.

992 **5.3.3. Expected Benefits**

993 A model solution with sequencers could facilitate MUD’s adoption within the bioeconomy and
994 significantly improve security across a large range of stakeholders including commercial,
995 academic, and government organizations.

996 By increasing security across multiple stakeholders, MUD would reduce the likelihood for
997 ransomware attacking and preventing usage of these valuable assets, as well as possible
998 intellectual property loss or privacy loss from exfiltration of data.

999 **5.4. Security Guidelines for Data Analysis Pipelines**

1000 Security guidelines (such as DISA STIGS or configuration baselines) manage and reduce
1001 cybersecurity risk by providing consensus-based and auditable configurations and security
1002 features for systems. They are a well-established method of providing transparency and
1003 alignment between user requirements and manufacturers’ product offerings. These guidelines

1004 may be used to harden data analysis pipelines and devices (for example, sequencers) to a trusted
1005 configuration baseline.

1006 **5.4.1. Use Case Description**

1007 Data analysis pipelines interact with a combination of open-source, off-the-shelf, and custom
1008 software including operating systems. Changes to software or configurations must be carefully
1009 validated and verified to ensure that they do not affect the pipelines' operation or introduce
1010 cybersecurity or privacy risks. Data analysis pipelines, including sequencers, typically prioritize
1011 ease of use and accuracy of results over minimizing cybersecurity and privacy risks. As such, it
1012 has been demonstrated that these pipelines and sequencers can be susceptible to adversarial
1013 threats [\[69\]](#).

1014 **5.4.2. Solution Idea**

1015 Security guidelines for data analysis pipelines (including sequencers) would provide best
1016 practice cybersecurity hardening guidance. These customized guidelines could be based off
1017 existing STIGs [\[70\]](#) or Center for Internet Security (CIS) Benchmarks [\[71\]](#) and tailored for the
1018 unique requirements of sequencers.

1019 These security guidelines might be modeled after NISTIR 8259 [\[72\]](#) for IoT devices that
1020 enumerates the capabilities, features, and functionalities that manufacturers of sequencers need
1021 to provide so that users can mitigate their cybersecurity risks. Security guidelines could include a
1022 requirement for the SBOM detailing all software installed on the sequencer to assist in
1023 identifying and mitigating cybersecurity vulnerabilities [\[58\]](#).

1024 **5.4.3. Expected Benefits**

1025 By establishing these guidelines and demonstrating feasibility, users would be able to include
1026 security standards in their purchasing requirements and sequencer manufacturers would have
1027 clear guidance on how to achieve the security standards. It would also enable assessing delivered
1028 products, enabling a device to be more quickly put into use both at initial delivery and after
1029 upgrade or service. This solution enables increased commerce by aligning user cybersecurity
1030 needs with manufacturers' cybersecurity offerings. The reduced friction from eliminating
1031 differing user expectations and manufacturer offerings would result in fewer purchasing delays
1032 and more competition between manufacturers on relevant user cybersecurity needs. It would also
1033 allow users to better manage their security risks in accordance with Federal information security
1034 requirements.

1035 **5.5. Demonstration Project for Genomic Data Risk Management**

1036 Because of the high value of genomic data, organizations in the bioeconomy and researchers
1037 using genomic data need to apply appropriate risk management. The NIST Risk Management
1038 Framework suite of guidance documents are comprehensive and can be applied to genomic data.

1039 A demonstration project providing an example of how they can be used along with the resources
1040 required for the effort would be beneficial.

1041 There does not exist a roadmap or demonstration project specifically tailored to the unique needs
1042 of aggregating and processing genomic data. These needs include consent management and
1043 reidentification risk. Consent management can be particularly complex in the large aggregations
1044 of subject data that are typical in oncology and precision medicine. The genomic data is typically
1045 pooled from many studies that may be collected with differing informed consent. The informed
1046 consent of each subject needs to be tracked and each analysis needs to be cleared for specific
1047 informed consents. The reidentification risk from genomic data stems from the characteristic that
1048 deidentification cannot be effectively done without sacrificing the utility of the data. If
1049 reidentification is successful, the subject data can be used for kinship, phenotype prediction,
1050 health information, and other possible future applications that may cause harm to the subject or
1051 the subject's kin.

1052 **5.5.1. Use Case Description**

1053 The target stakeholder would be a small organization, such as a university or start-up, that wants
1054 to aggregate and process genomic data. A start-up has important intellectual property that needs
1055 to be protected. A university may want to utilize NIH genomic data and must implement NIH
1056 cybersecurity and privacy requirements. Both need a solution to manage the risk of aggregating
1057 and processing genomic data that can be done in a timely and cost-effective manner.

1058 **5.5.2. Solution Idea**

1059 The demonstration project will identify controls, aligned to NIST SP 800-53, that are in place in
1060 a vendor-proposed solution that align to the unique requirements of handling genomic data. Two
1061 solutions exist that map their security implementation to NIST SP 800-53, [Terra.bio](#) and
1062 [DNAnexus](#). However, organizations must understand the unique requirements for handling
1063 genomic data. This does not alleviate the responsibility of the organization to implement
1064 appropriate risk management but can streamline the effort by using existing resources.

1065 **5.5.3. Expected Benefits**

1066 This demonstration project can identify gaps in vendor-proposed solutions aligned to genomic
1067 data handling cybersecurity and privacy requirements. The project could document savings in
1068 time, person hours, required expertise and documentation, to help an organization properly plan
1069 and budget their risk management effort. The guidance documentation produced from this effort
1070 would be specific to genomic data and could assist users by reducing mistakes due to missed
1071 requirements or misunderstanding of the requirements.

1072 **5.6. Demonstration Project for Analysis of Genomic Data Using Privacy** 1073 **Preserving and Enhancing Technologies**

1074 There are several promising technologies that may assist in addressing data subject privacy
1075 [\[3\]](#)[\[74\]](#)[\[75\]](#) and the breach of data confinement in genomic data sharing and analysis. Sharing
1076 high-value data typically requires a high degree of trust between data sharing organizations with
1077 significant contractual agreements. The process of negotiating and agreeing to the contractual
1078 arrangements and necessary controls can greatly slow and, in many cases, prohibit the sharing of
1079 data.

1080 Solutions that have been used or proposed for genomic datasets for addressing subject data
1081 privacy and preventing the breach of data confinement include:

- 1082 • A genomic analysis platform brings researchers to a secure location, requiring the use of
1083 plaintext and blocking all egress. This is not typically practical and severely limits
1084 research because all researchers must come to the secure location.
- 1085 • A “Secret Store” requires the use of predefined executables with no direct visibility to the
1086 data. However, this severely limits analysis options (only previously defined, validated,
1087 and verified executables can be used) available to researchers [\[4\]](#).
- 1088 • All data is stored so that analysis is done in a cloud environment created by a trusted
1089 authority, and all activity in that environment is monitored. This is not a complete
1090 solution, as researchers can provide visibility and/or their authorization to others.
1091 Additionally, this removes the advantage of remote access service which limits
1092 authorized users to only what they are authorized to do with the data, as currently
1093 available cloud solutions use identity-based authorization.
- 1094 • Federated Machine Learning (FML) [\[76\]](#), where collaborative learning is done among
1095 individual local nodes and shared with a centralized processor to produce an aggregate
1096 model, has been shown to be a promising method in precision medicine. It provides the
1097 capability to train models on competing or otherwise separate entities to produce an
1098 average model that captures the patterns present in the local nodes. However, unless
1099 combined with other privacy-enhancing technologies and controls, it can be vulnerable to
1100 data reconstruction attacks.²
- 1101 • Differential privacy can provide privacy guarantees and be combined with FML [\[77\]](#) or
1102 used to produce synthetic data [\[78\]](#) for machine learning. Differential private synthetic
1103 data adds noise to datasets so that an individual’s data cannot be distinguished within the
1104 dataset and allows the data to be analyzed and shared. Additionally, care must be taken
1105 when sharing this data, particularly if combining it with other data sets, as there is not yet
1106 sufficient research to determine whether combining it with other datasets maintains the
1107 privacy guarantee provided by differential privacy. It also allows for periodic updates as

² Privacy enhancing technologies or PETs include tools like homomorphic encryption, secure multi-party computation, and differential privacy. This is an evolving area and while some tools show promise, they are not a complete solution on their own and must be deployed with other controls. NIST is one of the organizers of the U.K.-U.S.PETs Prize Challenges which has a goal of “Accelerating the adoption and development of PETs.” Additional information and results are available at: <https://petsprizechallenges.com/>.

1108 the data evolves because differential private synthetic data can cope with the progress of
1109 the underlying dataset seamlessly. Researchers, however, have expressed concerns with
1110 the accuracy of the synthetic data as there is a tradeoff depending on where the privacy
1111 parameter is set. There are significant technology improvements in development to
1112 increase the accuracy of the synthetic data and some analyses are not as sensitive to the
1113 accuracy of the synthetic data.

- 1114 • Much progress has been made in privacy enhancing cryptography (PEC) [79] which
1115 addresses both data confinement and subject data privacy. Fully homomorphic encryption
1116 (FHE) allows computation on encrypted data and progress has been rapid at addressing
1117 its principal drawback of increased computational complexity. Another promising PEC
1118 technique is secure multi-party computation (SMPC) where analysis is performed over
1119 the data sets of several parties without revealing their input. Finally, there may be a place
1120 for zero-knowledge proofs (ZKP) [80][81] for validating results of genomic analysis
1121 without revealing the solution, which may fully preserve the privacy of the data
1122 underpinning the results.

1123 **5.6.1. Use Case Description**

1124 Large genomic datasets are required in oncology and precision medicine. Multiple dataset
1125 aggregation is often required. This typically requires the negotiation of multiple contracts
1126 between all the sharing organizations and the complexity of the contracts, since the datasets are
1127 extremely confidential, can be prohibitive. If data could be shared without a concern about
1128 exposure or exfiltration, the complexity of such contracts would be greatly simplified, and risk
1129 greatly reduced.

1130 **5.6.2. Solution Idea**

1131 Within the field of genomics, a solution that can virtually eliminate the risk of confidentiality or
1132 integrity loss of sharing genomic data between organizations and solve the confinement problem
1133 is federated multi-party homomorphic encryption [82]. This is a technique that allows for
1134 computation on encrypted data aggregated over multiple datasets. Exfiltration of the raw data is
1135 prevented as the authorized user is only able to obtain results from the computation which
1136 involves multiple datasets and authorized users cannot access the raw data in plaintext.

1137 At this stage in technology development, federated homomorphic encryption is not a general
1138 solution. It has been proven useful for several important problems in the analysis of genomic
1139 data, particularly in oncology and precision medicine research. In oncology, it has been shown to
1140 allow survival analysis that enables the effectiveness of different treatment options based on the
1141 genomic information of the patient and/or tumor to be compared [82]. In precision medicine, it
1142 allows genome-wide association studies (GWAS) [82] as well as machine learning (ML) training
1143 and testing. The homomorphic encryption solution is being used in Switzerland. A demonstration
1144 project in the United States working with The Broad Institute (United States) and/or Tune Insight
1145 (Switzerland) would be helpful for the technique to gain wider acceptance in the U.S.

1146 **5.6.3. Expected Benefits**

1147 This solution can enable more rapid progress in oncology treatment and precision medicine by
1148 allowing collaboration between organizations without complex contract negotiations or sharing
1149 agreements because the risk of exfiltration of data or privacy violation is virtually eliminated.
1150 Additionally, patient and research subjects' privacy would be controlled by technical means,
1151 which can be less prone to failure and data breaches compared to contractual controls.

1152 **6. Areas for Further Research**

1153 There are several problems where more research is needed that may be fruitful areas for NIST or
1154 NCCoE sponsored research work to address:

- 1155 • FedRAMP assessors use vulnerability scanners to assess the security posture of systems.
1156 Container (e.g., Docker or Pods) cybersecurity vulnerability assessment scanners are
1157 plagued with a high number of false positives. Research on how to improve the ability of
1158 vulnerability scanners to identify security issues in containers would benefit security
1159 professionals who maintain systems processing genomic data and optimize the use of
1160 their cybersecurity resources.
- 1161 • Research on how to securely integrate whole genomic data with a patient’s EHR while
1162 allowing for privacy and interoperability is needed. The FHIR standard provides a
1163 framework, but a genomic specific implementation is needed along with a demonstration
1164 project.
- 1165 • While federated multi-party homomorphic encryption can solve the confinement problem
1166 for a class of analyses, more work is needed as there are analysis methods not addressed
1167 by this solution. For the last few years, NIH has funded iDASH (integrating data for
1168 analysis, anonymization, and sharing) [\[74\]](#) and held secure genome analysis competition
1169 workshops to study specific cybersecurity mechanisms for genomic data systems. The
1170 key focus of this annual workshop is to advance cybersecurity technologies that are
1171 essential for genomic and biomedical system security and privacy. While such technical
1172 advancement is crucial for the security and privacy of genomic data sharing, more
1173 holistic and systematic security solutions that may use those technologies are also
1174 important for addressing systemic issues found in genomic information systems.

1175 **7. Conclusion**

1176 This document describes the challenges and concerns associated with handling genomic data, the
1177 current state of relevant cybersecurity and privacy risk management practices, gaps in
1178 implementing genomic data protections, and potential solutions along with areas for further
1179 research. Because of the value of genomic data, the associated challenges and concerns may
1180 affect national security, the U.S. economy, intellectual property, individual privacy, and under-
1181 protected populations. Existing cybersecurity and privacy risk management practices such as the
1182 NIST RMF, Cybersecurity Framework, and Privacy Framework must be tailored to effectively
1183 implement appropriate genomic data protections. Additionally, gaps persist in current policy,
1184 legislation, technology, and guidance for protecting genomic data.

1185 Solutions identified to address these challenges and gaps include:

- 1186 • A NIST Cybersecurity Framework Profile for handling of genomic data could help
1187 organizations identify gaps and prioritize investments in cybersecurity capabilities and
1188 controls.
- 1189 • A NIST Privacy Framework Profile for handling of genomic data could provide
1190 clarification on how to manage the privacy risks inherent in the aggregation, storage, and
1191 processing of human genomic data.
- 1192 • The Manufacturer Usage Description specification could improve sequencer security and
1193 reduce the likelihood of ransomware attacks as well as intellectual property loss or
1194 privacy loss from exfiltration of data.
- 1195 • Security guidelines or benchmarks for sequencers could provide best practice
1196 cybersecurity hardening guidance, require SBOMs to improve supply chain security, and
1197 improve cyber resiliency against future threats.
- 1198 • A demonstration project highlighting the benefits of using RMF guidance for protecting
1199 genomic data could illustrate how organizations can leverage appropriate secured cloud-
1200 based solutions to reduce the time, person hours, required expertise, and documentation
1201 required to implement effective cybersecurity and privacy practices.
- 1202 • A federated homomorphic encryption demonstration project for analysis of genomic data
1203 in precision medicine or oncology could illustrate how these solutions reduce the risk of
1204 confidentiality or integrity loss when sharing genomic data between organizations and
1205 help address the confinement problem.

1206 Future research areas may include methods for securely integrating whole genomic data with a
1207 patient's EHR while allowing for privacy and interoperability; improving the precision of
1208 vulnerability scanners for software containers; and technical solutions to solve the containment
1209 problem in genomic data for analysis methods not currently addressed by federated multi-party
1210 homomorphic encryption.

1211 **References**

- 1212 [1] Executive Office of the President, United States, "Executive order 14801: Advancing
1213 Biotechnology and Biomanufacturing Innovation for a Sustainable, Safe, and Secure
1214 American Bioeconomy," *Federal Register* 56849-56860, vol. 87, no. 178, 15 September
1215 2022.
- 1216 [2] National Institute of Standards and Technology (NIST), "NIST Privacy Framework, Version
1217 1.0," 16 January 2020. [Online]. Available: <https://doi.org/10.6028/NIST.CSWP.01162020>.
- 1218 [3] National Cybersecurity Center of Excellence, "NCCoE Virtual Workshop on the
1219 Cybersecurity of Genomic Data," 26 January 2022. [Online]. Available:
1220 [https://www.nccoe.nist.gov/get-involved/attend-events/nccoe-virtual-workshop-](https://www.nccoe.nist.gov/get-involved/attend-events/nccoe-virtual-workshop-cybersecurity-genomic-data/post-workshop-materials)
1221 [cybersecurity-genomic-data/post-workshop-materials](https://www.nccoe.nist.gov/get-involved/attend-events/nccoe-virtual-workshop-cybersecurity-genomic-data/post-workshop-materials).
- 1222 [4] National Cybersecurity Center of Excellence, "NCCoE Virtual Workshop on Exploring
1223 Solutions for the Cybersecurity of Genomic Data," 18 May 2022. [Online]. Available:
1224 [https://www.nccoe.nist.gov/get-involved/attend-events/nccoe-virtual-workshop-exploring-](https://www.nccoe.nist.gov/get-involved/attend-events/nccoe-virtual-workshop-exploring-solutions-cybersecurity-genomic-data)
1225 [solutions-cybersecurity-genomic-data](https://www.nccoe.nist.gov/get-involved/attend-events/nccoe-virtual-workshop-exploring-solutions-cybersecurity-genomic-data).
- 1226 [5] M. Gymrek and et.al., "Identifying Personal Genomes by Surname Inference," *Science*, vol.
1227 339, pp. 321-324, January 2013.
- 1228 [6] M. Naveed and et.al., "Privacy in the Genomic Era," *ACM Computer Survey*, vol. 48, no. 1,
1229 pp. 1-49, September 2015.
- 1230 [7] Y. Erlich and A. Narayanan, "Routes for breaching and protecting genetic privacy," *Nature*
1231 *Reviews*, vol. 15, pp. 409-421, June 2014.
- 1232 [8] Y. Erlich and et.al., "Identity Inference of Genomic Data Using Long-Range Familial
1233 Searches," *Science*, vol. 362, pp. 690-694, 2018.
- 1234 [9] C. A. Hutchison, "DNA sequencing: bench to bedside and beyond," 12 September 2007.
1235 [Online]. Available: <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC2094077/>. [Accessed
1236 15 December 2021].
- 1237 [10] National Institutes of Health (NIH), "DNA Sequencing Costs: Data," 1 November 2021.
1238 [Online]. Available: [https://www.genome.gov/about-genomics/fact-sheets/DNA-](https://www.genome.gov/about-genomics/fact-sheets/DNA-Sequencing-Costs-Data)
1239 [Sequencing-Costs-Data](https://www.genome.gov/about-genomics/fact-sheets/DNA-Sequencing-Costs-Data). [Accessed 15 December 2021].
- 1240 [11] NIH National Center for Biotechnology Information (NCBI), "Genome Reference
1241 Consortium," [Online]. Available: <https://www.ncbi.nlm.nih.gov/grc/human>. [Accessed 18
1242 July 2022].
- 1243 [12] J. Wagner and et.al., "Benchmarking challenging small variants with linked and long reads,"
1244 *Cell Genomics*, vol. 11, no. 2, p. 100128, 2022.

- 1245 [13] J. Chen and et.al., "Systematic comparison of germline variant calling pipelines cross
1246 multiple next-generation sequencers.," *Scientific Reports*, vol. 9, no. 1, pp. 1-13, 2019.
- 1247 [14] Z. Chen and et.al., "Systematic comparison of somatic variant calling performance among
1248 different sequencing depth and mutation frequency.," *Scientific Reports*, vol. 10, no. 1, pp.
1249 1-9, 2020.
- 1250 [15] Market Research Guru, "Global And United States Direct To Consumer Dtc Genetic
1251 Testing Market - Industry Reports," 20 May 2022 [Online]. Available:
1252 [https://www.marketresearchguru.com/global-and-united-states-direct-to-consumer-dtc-](https://www.marketresearchguru.com/global-and-united-states-direct-to-consumer-dtc-genetic-testing-market-20993717)
1253 [genetic-testing-market-20993717](https://www.marketresearchguru.com/global-and-united-states-direct-to-consumer-dtc-genetic-testing-market-20993717). [Accessed 20 July 2022].
- 1254 [16] NIH National Center for Biotechnology Information (NCBI), "Summary Statistics of dbGaP
1255 Data," [Online]. Available: [https://www.ncbi.nlm.nih.gov/projects/gap/summaries/cgi-](https://www.ncbi.nlm.nih.gov/projects/gap/summaries/cgi-bin/molecularDataPieSummary.cgi)
1256 [bin/molecularDataPieSummary.cgi](https://www.ncbi.nlm.nih.gov/projects/gap/summaries/cgi-bin/molecularDataPieSummary.cgi). [Accessed 20 July 2022].
- 1257 [17] M. O'Brien, "Who Has the Largest DNA Database? (2022)," 17 July 2021. [Online].
1258 Available: <https://www.dataminingdna.com/who-has-the-largest-dna-database/>. [Accessed
1259 15 July 2022].
- 1260 [18] T. Callaghan, "Responsible genetic genealogy," *Science*, vol. 366, no. 6462, p. 155, 2019.
- 1261 [19] R. Kain and et.al., "Database shares that transform research subjects into partners," *Nature*
1262 *biotechnology*, vol. 37, no. 10, pp. 1112-5, 2019.
- 1263 [20] K. Kozminski, "Biosecurity in the age of Big Data: a conversation with the FBI," *Molecular*
1264 *Biology of the Cell*, vol. 26, pp. 3894-3897, 2015.
- 1265 [21] National Counterintelligence and Security Center (NCSC), "China's Collection of Genomic
1266 and Other Healthcare Data From America: Risks to Privacy and U.S. Economic and
1267 National Security," pp. 1-5, February 2021.
- 1268 [22] A. Cohen, "MGH data breach exposes 10,000 patients," *Boston Herald*, 22 August 2019.
- 1269 [23] N. Grant, "DNA testing service Vitagene exposed thousands of customer records online for
1270 years," *Los Angeles Times*, 9 July 2019.
- 1271 [24] MyHeritage, "Security alert: malicious phishing attempt detected, possibly connected to
1272 GEDmatch breach," 21 July 2020. [Online]. Available:
1273 [https://blog.myheritage.com/2020/07/security-alert-malicious-phishing-attempt-detected-](https://blog.myheritage.com/2020/07/security-alert-malicious-phishing-attempt-detected-possibly-connected-to-gedmatch-breach/)
1274 [possibly-connected-to-gedmatch-breach/](https://blog.myheritage.com/2020/07/security-alert-malicious-phishing-attempt-detected-possibly-connected-to-gedmatch-breach/). [Accessed 16 November 2021].
- 1275 [25] "Devious 'Tardigrade' Malware Hits Biomanufacturing Facilities," 22 November 2021.
1276 [Online]. Available: <https://www.wired.com/story/tardigrade-malware-biomanufacturing/>.
1277 [Accessed 11 July 2022].

- 1278 [26] M. Foley, "Genetically Engineered Bioweapons: A New Breed of Weapons for Modern
1279 Warfare," *Dartmouth Undergraduate Journal of Science*, vol. Winter, 2013.
- 1280 [27] J. Peccoud and et.al., "Cyberbiosecurity: From Naive Trust to Risk Awareness," *Trends in*
1281 *Biotechnology*, vol. 36, no. 1, pp. 4-7, January 2018.
- 1282 [28] A. Schwab and et.al., "Genomic Privacy," *Clinical Chemistry*, vol. 64, no. 12, pp. 1696-
1283 1703, 2018.
- 1284 [29] M. Gianfrancesco and et.al., "Potential Biases in Machine Learning Algorithms Using
1285 Electronic Health Record Data," *JAMA Internal Medicine*, vol. 11, no. 178, pp. 1544-1547,
1286 2018.
- 1287 [30] Y. Joly and et.al., "Establishing the International Genetic Discrimination Observatory,"
1288 *Nature Genetics*, vol. 52, pp. 466-468, 2020.
- 1289 [31] R. Parikh and et.al., "Addressing Bias in Artificial Intelligence in Health Care," *JAMA*, vol.
1290 24, no. 322, pp. 2377-2378, 2019.
- 1291 [32] P. Grother and et.al., "NISTIR 8280 Face Recognition Vendor Test (FRVT) Part 3:
1292 Demographic Effects," December 2019. [Online]. Available:
1293 <https://nvlpubs.nist.gov/nistpubs/ir/2019/NIST.IR.8280.pdf>. [Accessed 28 July 2022].
- 1294 [33] CBS News, "The clones of polo," [Online]. Available: [https://www.cbsnews.com/news/the-](https://www.cbsnews.com/news/the-clones-of-polo)
1295 [clones-of-polo](https://www.cbsnews.com/news/the-clones-of-polo). [Accessed 20 July 2022].
- 1296 [34] National Institute of Health/National Human Genome Research Institute, "Eugenics and
1297 Scientific Racism," 18 May 2022. [Online]. Available: [https://www.genome.gov/about-](https://www.genome.gov/about-genomics/fact-sheets/Eugenics-and-Scientific-Racism)
1298 [genomics/fact-sheets/Eugenics-and-Scientific-Racism](https://www.genome.gov/about-genomics/fact-sheets/Eugenics-and-Scientific-Racism). [Accessed 19 July 2022].
- 1299 [35] National Institute of Standards and Technology (NIST), "SP 800-37 Rev. 2 Risk
1300 Management Framework for Information Systems and Organizations: A System Life Cycle
1301 Approach for Security and Privacy," December 2018. [Online]. Available:
1302 <https://csrc.nist.gov/publications/detail/sp/800-37/rev-2/final>.
- 1303 [36] Executive Office of the President, United States, "Improving the Nation's Cybersecurity," 12
1304 May 2021. [Online]. Available: [https://www.federalregister.gov/documents](https://www.federalregister.gov/documents/2021/05/17/2021-10460/improving-the-nations-cybersecurity)
1305 [/2021/05/17/2021-10460/improving-the-nations-cybersecurity](https://www.federalregister.gov/documents/2021/05/17/2021-10460/improving-the-nations-cybersecurity).
- 1306 [37] Global Alliance for Genomics and Health (GA4GH), "Framework for Responsible Sharing
1307 of Genomic and Health-Related Data," Global Alliance for Genomics and Health, 3
1308 September 2019. [Online]. Available: [https://www.ga4gh.org/genomic-data-](https://www.ga4gh.org/genomic-data-toolkit/regulatory-ethics-toolkit/framework-for-responsible-sharing-of-genomic-and-health-related-data/)
1309 [toolkit/regulatory-ethics-toolkit/framework-for-responsible-sharing-of-genomic-and-health-](https://www.ga4gh.org/genomic-data-toolkit/regulatory-ethics-toolkit/framework-for-responsible-sharing-of-genomic-and-health-related-data/)
1310 [related-data/](https://www.ga4gh.org/genomic-data-toolkit/regulatory-ethics-toolkit/framework-for-responsible-sharing-of-genomic-and-health-related-data/). [Accessed 15 November 2021].
- 1311 [38] U.S. Office of the Federal Register, "Public Law 113-282 - Federal Information Security
1312 Modernization Act of 2014," U.S. Government Publishing Office, 2014.

- 1313 [39] "FedRAMP," U.S. General Services Administration, [Online]. Available: [FedRAMP.gov](https://www.fedramp.gov).
1314 [Accessed 2023 01 18].
- 1315 [40] National Institute of Standards and Technology (NIST), "NIST SP 800-161r1," May 2022.
1316 [Online]. Available: <https://nvlpubs.nist.gov/nistpubs/SpecialPublications/NIST.SP.800-161r1.pdf>.
1317
- 1318 [41] Scott Rose and et.al., "NIST SP 800-27," August 2020. [Online]. Available:
1319 <https://csrc.nist.gov/publications/detail/sp/800-207/final>.
- 1320 [42] "CISA Zero Trust Maturity Model," Cybersecurity & Infrastructure Security Agency, 2022.
1321 [Online]. Available: <https://www.cisa.gov/zero-trust-maturity-model>.
- 1322 [43] OneTrust, "The Definitive Guide to Schrems II," 25 March 2022. [Online]. Available:
1323 <https://www.dataguidance.com/resource/definitive-guide-schrems-ii>. [Accessed 28 July
1324 2022].
- 1325 [44] E. Rafaelof and et.al., "Translation: China's 'Data Security Law (Draft)'," 2 June 2020.
1326 [Online]. Available: [https://www.newamerica.org/cybersecurity-
1327 initiative/digichina/blog/translation-chinas-data-security-law-draft/](https://www.newamerica.org/cybersecurity-initiative/digichina/blog/translation-chinas-data-security-law-draft/). [Accessed 2020 July
1328 2022].
- 1329 [45] N. Dhavate and M. Ramakant, "A look at proposed changes to India's (Personal) Data
1330 Protection Bill," 16 December 2021. [Online]. Available: [https://iapp.org/news/a/a-look-at-
1331 proposed-changes-to-indias-personal-data-protection-bill/](https://iapp.org/news/a/a-look-at-proposed-changes-to-indias-personal-data-protection-bill/). [Accessed 27 July 2022].
- 1332 [46] India Ministry of Electronics and Information Technology, "THE PERSONAL DATA
1333 PROTECTION BILL, Bill No. 373 of 2019," 2019. [Online]. Available:
1334 http://164.100/47/4/BillsTexts/LSBillTexts/Asintroduced/373_2019_LS_Eng.pdf.
1335 [Accessed 1 Aug 2022].
- 1336 [47] R. Creemers and et.al., "Translation: Cybersecurity Law of the People's Republic of China
1337 (Effective June 1, 2017)," New America, 29 June 2018. [Online]. Available:
1338 [https://www.newamerica.org/cybersecurity-initiative/digichina/blog/translation-
1339 cybersecurity-law-peoples-republic-china/](https://www.newamerica.org/cybersecurity-initiative/digichina/blog/translation-cybersecurity-law-peoples-republic-china/). [Accessed 1 August 2022].
- 1340 [48] D. Cyranoski, "China announces hefty fines for unauthorized collection of DNA," *Nature*
1341 *News*, 14 June 2019.
- 1342 [49] InCountry Staff, "Russian Data Protection Laws: Essential Guide on Compliance
1343 Requirements in Russia," InCountry, 19 March 2021. [Online]. Available:
1344 [https://incountry.com/blog/russian-data-protection-laws-essential-guide-on-compliance-
1345 requirements-in-russia/](https://incountry.com/blog/russian-data-protection-laws-essential-guide-on-compliance-requirements-in-russia/). [Accessed 1 August 2022].
- 1346 [50] OneTrust, "Russia - Data Protection Overview," OneTrust, April 2022. [Online]. Available:
1347 <https://www.dataguidance.com/notes/russia-data-protection-overview>. [Accessed 1 August
1348 2022].

- 1349 [51] Global Alliance for Genomics and Health (GA4GH), "Data Security Infrastructure Policy,"
1350 21 October 2019. [Online]. Available: [https://github.com/ga4gh/data-](https://github.com/ga4gh/data-security/blob/master/DSIP/DSIP_v4.0.md)
1351 [security/blob/master/DSIP/DSIP_v4.0.md](https://github.com/ga4gh/data-security/blob/master/DSIP/DSIP_v4.0.md).
- 1352 [52] National Institute of Standards and Technology (NIST), "Cybersecurity Framework Version
1353 1.1," 16 April 2018. [Online]. Available: <https://doi.org/10.6028/NIST.CSWP.04162018>.
1354 [Accessed March 2022].
- 1355 [53] National Institute of Standards and Technology, "NIST Special Publication 800-207 Zero
1356 Trust Architecture," 11 December 2020. [Online]. Available:
1357 <https://csrc.nist.gov/publications/detail/sp/800-207/final>.
- 1358 [54] "Security Technical Implementation Guides (STIGs)," DoD Cyber Exchange Public,
1359 [Online]. Available: <https://public.cyber.mil/stigs/>. [Accessed 18 January 2023].
- 1360 [55] "bioisac," 2022. [Online]. Available: <https://www.isac.bio/>. [Accessed 18 January 2023].
- 1361 [56] Future of Privacy Forum, "Privacy Best Practices for Consumer Genetic Testing Services,"
1362 31 July 2018. [Online]. Available: [https://fpf.org/wp-content/uploads/2018/07/Privacy-Best-](https://fpf.org/wp-content/uploads/2018/07/Privacy-Best-Practices-for-Consumer-Genetic-Testing-Services-FINAL.pdf)
1363 [Practices-for-Consumer-Genetic-Testing-Services-FINAL.pdf](https://fpf.org/wp-content/uploads/2018/07/Privacy-Best-Practices-for-Consumer-Genetic-Testing-Services-FINAL.pdf). [Accessed 2 December
1364 2021].
- 1365 [57] Global Alliance for Genomics and Health (GA4GH), "Privacy and Security Policy," 26 May
1366 2015. [Online]. Available: [https://www.ga4gh.org/wp-content/uploads/Privacy-and-](https://www.ga4gh.org/wp-content/uploads/Privacy-and-Security-Policy.pdf)
1367 [Security-Policy.pdf](https://www.ga4gh.org/wp-content/uploads/Privacy-and-Security-Policy.pdf).
- 1368 [58] Cybersecurity & Infrastructure Security Agency (CISA), "Software Bill of Materials,"
1369 [Online]. Available: <https://www.cisa.gov/sbom>. [Accessed 10 June 2022].
- 1370 [59] 110th U.S. Congress, "Public Law 110-233," 21 May 2008. [Online]. Available:
1371 <https://www.govinfo.gov/content/pkg/PLAW-110publ233/pdf/PLAW-110publ233.pdf>.
1372 [Accessed 2 December 2021].
- 1373 [60] 104th U.S. Congress, "Health Insurance Portability and Accountability Act of 1996," 20
1374 August 1996. [Online]. Available: [https://aspe.hhs.gov/reports/health-insurance-portability-](https://aspe.hhs.gov/reports/health-insurance-portability-accountability-act-1996)
1375 [accountability-act-1996](https://aspe.hhs.gov/reports/health-insurance-portability-accountability-act-1996). [Accessed 2 December 2021].
- 1376 [61] Department of Health and Human Services (HHS), "Guidance Regarding Methods for De-
1377 identification of Protected Health Information in Accordance with the Health Insurance
1378 Portability and Accountability Act (HIPAA) Privacy Rule," [Online]. Available:
1379 [https://www.hhs.gov/hipaa/for-professionals/privacy/special-topics/de-](https://www.hhs.gov/hipaa/for-professionals/privacy/special-topics/de-identification/index.html#standard)
1380 [identification/index.html#standard](https://www.hhs.gov/hipaa/for-professionals/privacy/special-topics/de-identification/index.html#standard). [Accessed 10 June 2022].
- 1381 [62] Department of Health and Human Services (HHS), "Federal Policy for the Protection of
1382 Human Subjects ('Common Rule')," 13 December 2022. [Online]. Available:
1383 <https://www.hhs.gov/ohrp/regulations-and-policy/regulations/common-rule/index.html>.
1384 [Accessed 18 January 2023].

- 1385 [63] Arizona State Legislature, "44-8002. Direct-to-consumer genetic testing company
1386 requirements; prohibition," 29 September 2021. [Online]. Available:
1387 <https://www.azleg.gov/viewdocument/?docName=https://www.azleg.gov/ars/44/08002.htm>.
1388 [Accessed 18 January 2023].
- 1389 [64] California State Legislature, "Senate Bill No. 41, Chapter 596," 7 October 2021. [Online].
1390 Available:
1391 https://leginfo.legislature.ca.gov/faces/billTextClient.xhtml?bill_id=202120220SB41.
1392 [Accessed 18 January 2023].
- 1393 [65] Utah Legislative General Counsel, "Genetic Information Privacy Act," 22 February 2021.
1394 [Online]. Available: <https://le.utah.gov/~2021/bills/sbillint/SB0227.pdf>. [Accessed 18
1395 January 2023].
- 1396 [66] National Institute of Standards and Technology (NIST), "Cybersecurity Framework -
1397 Examples of Framework Profiles," NIST, 5 July 2022. [Online]. Available:
1398 <https://www.nist.gov/cyberframework/examples-framework-profiles>. [Accessed 7 July
1399 2022].
- 1400 [67] N. Homer and et.al., "Resolving Individuals Contributing Trace Amounts of DNA to Highly
1401 Complex Mixtures Using High-Density SNP Genotyping Microarrays," *PLoS Genetics*, vol.
1402 4, no. 8, p. e1000167, 2007.
- 1403 [68] National Institute of Standards and Technology (NIST), "NIST Special Publication 1800-15
1404 Securing Small-Business and Home Internet of Things (IoT) Devices: Mitigating Network-
1405 Based Attacks Using Manufacturer Usage Description (MUD)," May 2021. [Online].
1406 Available: <https://nvlpubs.nist.gov/nistpubs/SpecialPublications/NIST.SP.1800-15.pdf>.
- 1407 [69] C. Cimpanu, "Mysterious Iranian group is hacking into DNA sequencers," 14 June 2019.
1408 [Online]. Available: [https://www.zdnet.com/article/mysterious-iranian-group-is-hacking-
1409 into-dna-sequencers/](https://www.zdnet.com/article/mysterious-iranian-group-is-hacking-into-dna-sequencers/). [Accessed 13 June 2022].
- 1410 [70] DoD Cyber Exchange Public, "STIGs Document Library," [Online]. Available:
1411 <https://public.cyber.mil/stigs/downloads/>. [Accessed 10 June 2022].
- 1412 [71] Center for Internet Security, "CIS Benchmarks," [Online]. Available:
1413 <https://www.cisecurity.org/cis-benchmarks/>. [Accessed 10 June 2022].
- 1414 [72] National Institute of Standards and Technology (NIST), "Foundational Cybersecurity
1415 Activities for IoT Device Manufacturers," May 2020. [Online]. Available:
1416 <https://csrc.nist.gov/publications/detail/nistir/8259/final>.
- 1417 [73] National Institute of Standards and Technology (NIST), "Cyber-Physical Systems/Internet
1418 of Things Testbed," 7 April 2022. [Online]. Available: [https://www.nist.gov/programs-
1419 projects/cyber-physical-systemsinternet-things-testbed](https://www.nist.gov/programs-projects/cyber-physical-systemsinternet-things-testbed). [Accessed 15 July 2022].

- 1420 [74] "iDASH Privacy & Security Workshop," [Online]. Available:
1421 <http://www.humangenomeprivacy.org>. [Accessed 19 November 2021].
- 1422 [75] X. Shi and X. Wu, "An overview of human genetic privacy," *Annals of the New York*
1423 *Academy of Sciences*, vol. 1387, no. 1, pp. 61-72, 2017.
- 1424 [76] N. Rieke and et.al., "The future of digital health with federated learning," *npj Digital*
1425 *Medicine*, vol. 3, no. 119, 2020.
- 1426 [77] M. A. Aziz and et.al., "Generalized Genomic Data Sharing for Differentially Private
1427 Federated Learning," *Journal of Biomedical Informatics*, no.
1428 <https://doi.org/10.1016/j.jbi.2022.104113>, 2022.
- 1429 [78] J. Near and D. Darais, "Differentially Private Synthetic Data," 3 May 2021. [Online].
1430 Available: [https://www.nist.gov/blogs/cybersecurity-insights/differentially-private-](https://www.nist.gov/blogs/cybersecurity-insights/differentially-private-synthetic-data)
1431 [synthetic-data](https://www.nist.gov/blogs/cybersecurity-insights/differentially-private-synthetic-data). [Accessed 28 July 2022].
- 1432 [79] National Institute of Standards and Technology (NIST), "Privacy-Enhancing Cryptography
1433 PEC," 22 July 2022. [Online]. Available: <https://csrc.nist.gov/projects/pec>. [Accessed 28
1434 July 2022].
- 1435 [80] S. Goldwasser and et.al., "The Knowledge Complexity of Interactive Proof Systems," *Siam*
1436 *Journal of Computing*, vol. 18, no. 1, pp. 186-208, 1989.
- 1437 [81] L. Lesavre and et.al., "Blockchain Networks: Token Design and Managment Overview,"
1438 National Institutue of Standards and Technology, NISTIR 8301, 2021.
- 1439 [82] D. Froelicher and et.al., "Truly privacy-preserving federated analytics for precision
1440 1441 medicine with multiparty homomorphic encryption.," *Nature Communications*, vol. 12,
1441 no. 1442 5910, 2021.

1442 **Appendix A. List of Symbols, Abbreviations, and Acronyms**

1443 The following acronyms are used in this publication.

1444 **AI**
1445 Artificial Intelligence

1446 **BIO-ISAC**
1447 Bioeconomy Information Sharing and Analysis Center

1448 **CIS**
1449 Center for Internet Security

1450 **CISA**
1451 Cybersecurity and Infrastructure Security Agency

1452 **CRISPR**
1453 Clustered Regularly Interspaced Short Palindromic Repeats

1454 **CRISPR-Cas**
1455 CRISPR-Associated Protein

1456 **DbGaP**
1457 Database of Genotypes and Phenotypes

1458 **DISA**
1459 Defense Information Systems Agency

1460 **DNA**
1461 Deoxyribonucleic acid

1462 **DTC**
1463 Direct-To-Consumer

1464 **EEA**
1465 European Economic Area

1466 **EHR**
1467 Electric Health Record

1468 **EO**
1469 Executive Order

1470 **FBI**
1471 Federal Bureau Investigation

1472 **FedRAMP**
1473 Federal Risk and Authorization Management Program

1474 **FHE**
1475 Fully Homomorphic Encryption

1476 **FHIR**
1477 Fast Healthcare Interoperability Resources

1478 **FIPPs**
1479 Fair Information Practice Principles

- 1480 **FISMA 2014**
- 1481 Federal Information Security Modernization Act (2002)

- 1482 **FISMA 2002**
- 1483 Federal Information Security Management Act (2002)

- 1484 **FPF**
- 1485 Future of Privacy Forum

- 1486 **GA4GH**
- 1487 Global Alliance for Genomic Health

- 1488 **GINA**
- 1489 Genetic Information Nondiscrimination Act (2008)

- 1490 **GWAS**
- 1491 Genome-Wide Association Studies

- 1492 **HIMSS**
- 1493 Healthcare Information and Management Systems Society

- 1494 **HIPAA**
- 1495 Health Insurance Portability and Accountability Act (1996; amended 2013)

- 1496 **iDASH**
- 1497 Integrating Data for Analysis, Anonymization, and Sharing

- 1498 **IoT**
- 1499 Internet of Things

- 1500 **IP**
- 1501 Intellectual Property

- 1502 **MUD**
- 1503 Manufacturer Usage Description

- 1504 **NCBC**
- 1505 National Centers for Biomedical Computing

- 1506 **NCBI**
- 1507 National Center for Biotechnology Information

- 1508 **NCCoE**
- 1509 NIST National Cybersecurity Center of Excellence

- 1510 **NGS**
- 1511 Next-Generation Sequencing

- 1512 **NISTIR**
- 1513 NIST Internal Report

- 1514 **OMB**
- 1515 Office of Management and Budget

- 1516 **PEC**
- 1517 Privacy Enhancing Cryptography

1518	PHI
1519	Protected Health Information
1520	PMO
1521	Program Management Office
1522	RMF
1523	Risk Management Framework
1524	RNA
1525	Ribonucleic Acid
1526	SBOM
1527	Software Bill of Materials
1528	SMPC
1529	Secure Multi-Party Computation
1530	SRA
1531	Sequence Read Archive
1532	STIGS
1533	Security Technical Implementation Guides
1534	The Common Rule
1535	Federal Policy for the Protection of Human Subjects
1536	VA
1537	Department of Veterans Affairs
1538	ZKP
1539	Zero-Knowledge Proof
1540	ZTA
1541	Zero-Trust Architecture