
Data Stewardship and Curation Practices in AI-Driven Genomics and Automated Microscopy Image Analysis for High-Throughput Screening Studies: Promoting Robust and Ethical AI Applications

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 The increasing adoption of AI and next-generation sequencing (NGS) has revolutionized
2 genomics and high-throughput screening (HTS), transforming how cellular
3 processes and disease mechanisms are understood. However, these advancements
4 generate vast datasets requiring effective data stewardship and curation practices to
5 maintain data integrity, privacy, and accessibility. This review consolidates exist-
6 ing knowledge on key aspects, including data governance, quality management,
7 privacy measures, ownership, access control, accountability, traceability, curation
8 frameworks, and storage systems. Major challenges such as managing biases, en-
9 suring data quality, and securing privacy are highlighted. Advanced cryptographic
10 techniques, federated learning, and blockchain technology are proposed as strategic
11 solutions, emphasizing standards compliance, ethical oversight, and tailored access
12 control frameworks. Effective data stewardship is vital for advancing AI-driven
13 genomics and microscopy research. Stakeholders must prioritize robust data gover-
14 nance and privacy measures to ensure data integrity and ethical use. Collaborative
15 efforts should focus on developing transparent data-sharing policies and interoper-
16 able platforms to foster innovation and advance research practices. The study
17 promotes collaboration among researchers, robust data governance, privacy and
18 security, clear policies, and educational initiatives to prepare future researchers.

19 1 Introduction

20 Advancements in AI and next-generation sequencing (NGS) have revolutionized
21 genomics and high-throughput screening (HTS) studies, enabling the integration of multi-dimensional data [6, 4, 5].
22 Automated high-content screening (HCS) methodologies, combining microscopy image acquisition
23 and analysis, are now pivotal for understanding cellular processes and assessing drug efficacy [3, 2].
24 However, these technologies generate vast datasets that require robust data stewardship and curation
25 practices to ensure their reliability and accessibility. Therefore, this study aimed to elucidate best
26 practices in data stewardship and curation for AI-driven genomics and automated microscopy image
27 analysis within high-throughput screening studies.

28 2 Methods

29 A systematic literature search was conducted up to December 30, 2023, across PubMed, MEDLINE,
30 Embase, Scopus, and Web of Science. The search focused on data governance, curation frameworks,

31 algorithmic bias, and data storage. Realist synthesis methodology was used to integrate diverse
 32 theoretical frameworks, with three independent reviewers. The review process included six stages,
 33 starting with an extensive search across multiple research databases, resulting in 273 documents. This
 34 was followed by screening based on broad criteria, titles, abstracts, and full texts, which narrowed the
 35 pool to 38 highly relevant citations.

36 **3 Experimental Results**

37 Our findings revealed a significant surge in research activity in recent years, particularly in 2023,
 38 reflecting the increasing recognition of the importance of robust data governance frameworks. Notably,
 39 while 36 articles extensively discussed data interoperability and sharing measures, areas such as
 40 model explainability and data augmentation remained underexplored, highlighting crucial gaps that
 41 need to be addressed. The integration of diverse data types, including sequencing, clinical, proteomic,
 42 and imaging data, underscored the complexity and breadth of AI applications in genomics and
 microscopy.

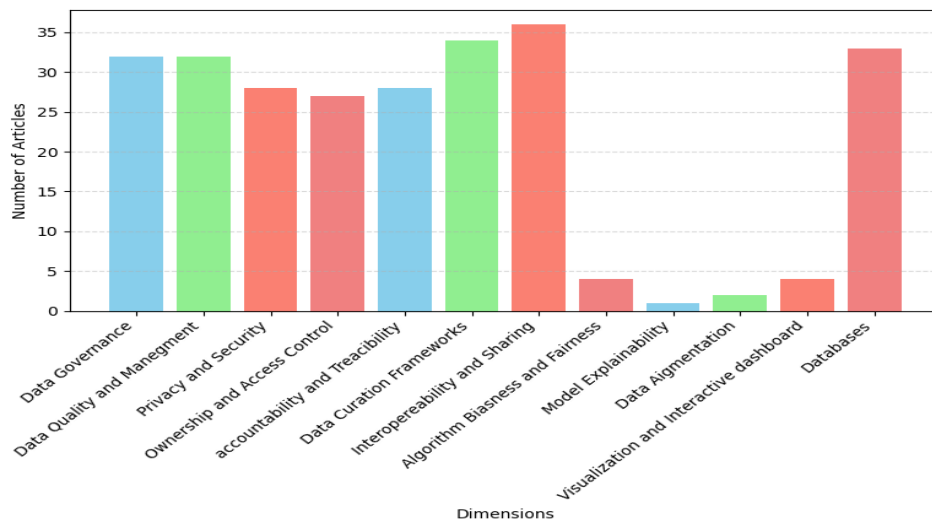


Figure 1: Illustrated the number of articles addressed the data stewardship and curation practice through different dimensions

43

44 Moreover, our review emphasized that infrastructure optimization, ethical considerations, access
 45 control mechanisms, and transparent data sharing policies were the most critical challenges in AI-
 46 based data stewardship. Advanced cryptographic techniques, federated learning, and blockchain
 47 technology were proposed to address challenges like data quality, privacy, and bias management. We
 48 identified robust data governance measures, such as GA4GH standards[3], DUO versioning, and
 49 attribute-based access control[2], as essential for ensuring data integrity, security, and ethical use.
 50 The importance of Data Management Plans (DMPs) [1], meticulous metadata curation, and advanced
 51 cryptographic techniques emerged as pivotal in mitigating data security and identifiability risks.

52 **4 Conclusion**

53 These findings provide a comprehensive overview of current practices and challenges in data stew-
 54 ardship, offering a roadmap for enhancing the robustness and ethical standards of AI applications in
 55 genomics and microscopy. Effective data stewardship and curation are vital for advancing AI-driven
 56 genomics and microscopy image analysis. Prioritizing robust governance, quality management, and
 57 secure sharing frameworks is essential. Collaborative efforts must focus on developing transparent
 58 data-sharing policies and interoperable platforms to foster innovation and advance research practices.

59 **References**

- 60 [1] Faisal M Fadlelmola, Lyndon Zass, Melek Chaouch, Chaimae Samtal, Verena Ras, Judit Ku-
61 muthini, Sumir Panji, and Nicola Mulder. Data management plans in the genomics research
62 revolution of africa: Challenges and recommendations. *Journal of biomedical informatics*,
63 122:103900, 2021.
- 64 [2] David Reddick, Justin Presley, Frank Alex Feltus, and Susmit Shannigrahi. Wip: Aabac-
65 automated attribute based access control for genomics data. In *Proceedings of the 27th ACM on*
66 *Symposium on Access Control Models and Technologies*, pages 217–222, 2022.
- 67 [3] Michael C Schatz, Anthony A Philippakis, Enis Afgan, Eric Banks, Vincent J Carey, Robert J
68 Carroll, Alessandro Culotti, Kyle Ellrott, Jeremy Goecks, Robert L Grossman, et al. Inverting
69 the model of genomics data sharing with the nhgri genomic data science analysis, visualization,
70 and informatics lab-space. *Cell Genomics*, 2(1), 2022.
- 71 [4] Marcel P Schilling, Razan El Khaled El Faraj, Joaquín Eduardo Urrutia Gómez, Steffen J
72 Sonnentag, Fei Wang, Britta Nestler, Véronique Orian-Rousseau, Anna A Popova, Pavel A
73 Levkin, and Markus Reischl. Automated high-throughput image processing as part of the
74 screening platform for personalized oncology. *Scientific Reports*, 13(1):5107, 2023.
- 75 [5] Rachel H Toczydlowski, Libby Liggins, Michelle R Gaither, Tanner J Anderson, Randi L Barton,
76 Justin T Berg, Sofia G Beskid, Beth Davis, Alonso Delgado, Emily Farrell, et al. Poor data
77 stewardship will hinder global genetic diversity surveillance. *Proceedings of the National*
78 *Academy of Sciences*, 118(34):e2107934118, 2021.
- 79 [6] Galen EB Wright, Pieter GJ Koornhof, Adebowale A Adeyemo, and Nicki Tiffin. Ethical and
80 legal implications of whole genome and whole exome sequencing in african populations. *BMC*
81 *Medical Ethics*, 14:1–15, 2013.