Databricks vs. Snowflake for Life Sciences: A Comparison

By InuitionLabs.ai · 10/17/2025 · 35 min read

databricks snowflake life sciences data lakehouse cloud data platform genomics data ai/ml clinical analytics bioinformatics

Sometimes of the science of t



Executive Summary

The life sciences industry is experiencing an unprecedented data deluge: genomic, clinical, and real-world healthcare data volumes are growing exponentially. For example, genomic data in public repositories jumped from ~47 GB in 2007 to ~28 PB by 2024 (a 620,000-fold increase) (pmc.ncbi.nlm.nih.gov), and healthcare data volumes are projected to grow ~36% per year in the next five years (medium.com). In this environment, modern cloud data platforms have become critical. Two leading solutions are **Databricks** (a cloud *data + Al lakehouse*) and **Snowflake** (a cloud *data warehouse/Cloud Data Platform*). Both platforms aim to enable advanced analytics, Al/ML, and data sharing, but they employ different architectures and strengths.

This report provides a deep, evidence-based comparison of Databricks vs. Snowflake in the context of life sciences. We cover their technological architectures, capabilities, and how each addresses life-science-specific needs (e.g. genomics data processing, clinical analytics, regulatory compliance, collaboration). We begin with background on data growth and the origins of the platforms, then examine detailed aspects including data types, compute, machine learning support, data sharing, and compliance features. We review case studies and industry examples (e.g. Pfizer, Regeneron, Anthem, Novartis) to illustrate real-world usage and outcomes. Datadriven insights and expert analyses are weaved throughout, highlighting metrics where available. Finally, we discuss future trends (such as Al integration and open data ecosystems) and offer guidance on choosing or combining these platforms for life sciences.

Key Findings: Databricks' Lakehouse unifies data engineering, data science, and ML on massive multi-modal datasets, making it well-suited for big-data R&D tasks like genomics, imaging, and Al model training (www.databricks.com) (www.prnewswire.com). Snowflake's Cloud Data Platform emphasizes high-performance SQL analytics, data sharing, and ease of use for structured data workloads, with built-in governance and security to handle sensitive patient and pharma data (venturebeat.com) (www.prnewswire.com). In practice, leading life sciences organizations often use both: for example, Snowflake may be leveraged for scalable data warehousing and collaboration (real-world evidence, biostatistics) while Databricks powers large-scale modelling and AI experiments (medium.com) (www.prnewswire.com). Both platforms now aggressively support AI/ML: Databricks via deep integration of Apache Spark, MLflow, and acquisitions (e.g. Dolly/MosaicML) (www.prnewswire.com) (techcrunch.com), and Snowflake via Snowpark, H2O, Streamlit/Neeva and partnerships (NVIDIA/NeMo) to enable model training/inference on cloud data (venturebeat.com) (techcrunch.com). Multiple case studies show significant benefits (e.g. Pfizer using Snowflake saw a 57% TCO reduction (www.snowflake.com); Databricks customers report much faster insights and eliminated data silos (www.prnewswire.com) (www.prnewswire.com)). Future trends such as generative AI, knowledge graphs, and federated data networks will further shape both platforms' roles in life sciences.

Introduction and Background

The Data Challenge in Life Sciences

Life sciences (including pharmaceuticals, biotechnology, genomics, and healthcare organizations) are data-intensive by nature. Advances in high-throughput sequencing, imaging, EHR systems, and IoT are generating **massive and diverse datasets**. Research literature growth underscores this: for example, PubMed saw ~7.3 million new articles from 2004–2013 (a 49% increase) (www.databricks.com). Next-generation sequencing alone has driven the cost of sequencing a human genome below \$1,000 and its data output to petabyte scales (pmc.ncbi.nlm.nih.gov). These trends mean scientists now integrate genomics, transcriptomics, proteomics, phenotypic, and clinical data to drive discovery. However, legacy platforms (on-premises databases or siloed data lakes) often fragment data and impede analytics.

Cloud platforms promise scalability, flexibility, and collaboration. Cloud vendors offer managed storage, high-performance compute, and pay-as-you-go pricing. Analysts note that cloud infrastructure can easily scale HPC resources (GPUs/TPUs) for ML while handling compliance (e.g. HIPAA, GDPR) via built-in security (pmc.ncbi.nlm.nih.gov) (venturebeat.com). For example, one study highlights that cloud computing provides "dynamic resource scaling" and "robust security with compliance (HIPAA, GDPR)" compared to on-premises systems (pmc.ncbi.nlm.nih.gov). Figure 1 (below) compares traditional on-premises/supercomputer environments with cloud computing, illustrating benefits of cloud for life science data (from Koreeda et al. 2024 (pmc.ncbi.nlm.nih.gov)).

Criteria	On-Premises	Supercomputers	Cloud Computing
Customization	Full hardware/software customization	Optimized for specific tasks	Dynamic resource scaling
Cost	High fixed costs (maintenance, HW)	Shared usage constraints, admin overhead	Pay-as-you-go (can reduce short-term costs) (pmc.ncbi.nlm.nih.gov)
Performance	Fast if invested; limited by HW	High for large-scale jobs	High with scalable virtual clusters (pmc.ncbi.nlm.nih.gov)
Flexibility	Limited by physical resources	Limited by batch scheduling	Scale up/down on demand (elastic) (pmc.ncbi.nlm.nih.gov)
Security	Managed internally (risky if misconfig)	Sensitive data handling challenging	Built-in compliance (HIPAA, etc.) (pmc.ncbi.nlm.nih.gov)
Data Sharing	Limited (often manual)	Hard for long-term storage sharing	Global sharing/collaboration (cloud objects) (pmc.ncbi.nlm.nih.gov)
Risk	High maintenance burden	Limited access windows	Vendor lock-in risk but lower ops burden (pmc.ncbi.nlm.nih.gov)

Figure 1: Comparison of on-premises, supercomputing, and cloud computing for large-scale biological data (adapted from Koreeda et al. 2024 (pmc.ncbi.nlm.nih.gov)).

Cloud adoption in life sciences is accelerating due to these benefits. For example, healthcare data volume is projected to grow ~36% annually over the next five years - faster than any other major industry (medium.com) - which drives the need for elastic cloud platforms. Compliance is also crucial: U.S. healthcare and pharma must meet HIPAA, GxP, SOC2, and other standards. Both Databricks and Snowflake explicitly support these needs: Snowflake's platform is built for security and compliance (venturebeat.com), and Databricks provides HIPAA-ready configurations (docs.databricks.com) and regulatory-grade MLOps (www.databricks.com).

Databricks and Snowflake: Origins and Architectures

Databricks was founded in 2013 by the creators of Apache Spark (from UC Berkeley) (dnamic.ai). It introduced the Data Lakehouse paradigm, blending data lakes and warehouses. The Databricks Lakehouse natively integrates data engineering and AI/ML: it is essentially a Spark-based unified analytics environment that can handle structured, semi-structured, and unstructured data at scale (using its Delta Lake format for transactional reliability). The platform is multi-cloud (runs on AWS, Azure, GCP) and supports Python, SQL, R, Scala and more (www.databricks.com) (www.databricks.com). Databricks emphasizes flexibility for data science workflows (e.g. via collaborative notebooks) and advanced ML.

Snowflake, founded in 2012, pioneered a true cloud data warehouse architecture (dnamic.ai). Snowflake is a SaaS platform built on top of cloud object storage (originally AWS S3, now also Azure Blob and Google Cloud Storage) with a unique multi-cluster, shared-data architecture. It separates compute from storage, enabling virtually unlimited concurrency via multi-cluster "virtual warehouses" and automatic scaling. Snowflake initially focused on structured data (highperformance SQL analytics) but later added support for semi-structured data (VARIANT type for JSON) and integrations for machine learning (Snowpark, Python UDFs, built-in ML and Data Science tools). It also offers a data exchange for secure sharing of data across organizations. Snowflake is multi-region and multi-cloud, allowing organizations to run a single platform across clouds. Snowflake's design prioritizes simplicity for BI workloads, strong performance on SQL queries, and enterprise-grade security/compliance (dnamic.ai) (venturebeat.com).

Both companies now pitch themselves as "data+AI" platforms for enterprises. Databricks markets its Lakehouse specifically for healthcare and life sciences: it claims to "consolidate massive volumes of data" across the drug development lifecycle, unlocking insights and lowering costs (www.databricks.com). Snowflake has likewise created a dedicated "Healthcare & Life Sciences Data Cloud" – a cross-cloud data platform tailored to the industry's challenges (data silos, compliance, need for collaboration) (venturebeat.com) (venturebeat.com). The concurrently rising momentum of both platforms has led some analysts to note the "data dogfight" between them; both aim to be one-stop-shops for enterprise data and Al (venturebeat.com).

Technical Comparison

Data Architecture and Storage

Databricks (Lakehouse): Databricks implements the *lakehouse* architecture, which layers a unified storage (often an object store or "Data Lake") and computes engine (Apache Spark™) into one platform (www.prnewswire.com). All data − structured tables and unstructured data (e.g. genomic files, medical images, text documents) − can reside together in cloud storage (AWS S3, Azure Data Lake, etc.) and be accessed with ACID guarantees via Delta Lake. Databricks emphasizes openness: it builds on open-source components (Spark, Delta Lake, MLflow) and supports arbitrary data types. This makes it well-suited for ingesting raw experimental data (FASTQ/BAM genomic files, pathology slide images, time-series sensor data) as well as tabular clinical data, all within a single "health lakehouse" (www.databricks.com) (www.prnewswire.com). With Databricks, life science teams can apply Spark-based processing pipelines (ETL, ML pipelines) on this integrated store, leveraging distributed computation and strong caching. Databricks also recently introduced features like Unity Catalog for data governance, and high-performance delta sharing for collaboration, though Snowflake leads in native cross-company sharing.

Snowflake (Data Cloud): Snowflake's data architecture is a managed multi-cluster service on the cloud. Data is stored in Snowflake's proprietary internal format on top of cloud storage. Unlike a traditional warehouse, it can directly query semi-structured formats (JSON, Avro) using SQL. Snowflake's key selling point is separation of compute and storage: storage autoscaling is handled independently, and compute clusters (virtual warehouses) can be spun up/down on demand or in parallel. Snowflake's architecture excels at *structured data warehousing*: large data sets (patient records, standardized clinical and experimental tabular data) can be loaded into Snowflake's tables and queried with high concurrency. Tables in Snowflake are automatically clustered and indexed by the system (via micro-partitions). Snowflake also offers features like Time Travel and Zero-Copy Cloning for data recovery and versioning, which we discuss later. Snowflake natively implements strong governance: all data is encrypted at rest and in flight, with robust role-based access controls (venturebeat.com).

In summary, **Databricks** provides a more flexible multi-modal lakehouse (strong for ingesting new data types, heavy processing), whereas **Snowflake** offers a streamlined SQL-based data cloud (strong for high-performance analytics on structured data). These differences emerge in practice: Databricks documentation highlights connecting "EHRs, wearables, imaging platforms, genome sequencers and more" into one view (www.databricks.com). Snowflake emphasizes its ability to "centralize, integrate, and exchange critical and sensitive data at scale" across clouds (venturebeat.com) (venturebeat.com). (In Tables 2 and 3 below, we compare key attributes of both platforms.)



Attribute	Databricks (Lakehouse)	Snowflake (Data Cloud)
Architecture	Data Lake + Delta Lake + Apache Spark engine (Lakehouse) (www.prnewswire.com) (www.databricks.com). All data (structured/unstructured) in unified storage; ACID via Delta.	Cloud Data Warehouse; one database across clouds. Separation of storage/compute; multicluster shared data architecture (dnamic.ai) (venturebeat.com).
Primary Data Models	Supports structured tables, semi-structured (JSON), unstructured (files, images) (www.databricks.com). Schema-on-read for raw data; Delta Tables (schema) for curated data.	Optimized for structured and semi- structured data (VARIANT JSON); can query files via external tables. Relational schema with optional semi-structured columns.
Compute Engine	Apache Spark core (distributed CPU), optimized for big data jobs. Also Databricks SQL (presto-like) and Photon (C++) (www.prnewswire.com). Autoscaling Spark clusters, GPU support for ML.	Native SQL engine (Mpp), automatically scales up by adding compute clusters. Snowpark for Python/Java UDFs (CPU, GPU in Partnered edge). Recent GPU/container support.
ML/AI Capabilities	Integrated MLflow, TensorFlow/PyTorch support in notebooks, ML pipelines. Strong in advanced analytics; supports real-time streaming analytics (www.databricks.com) (techcrunch.com).	Snowpark ML (Python, R) and built-in UDFs; partnerships (e.g. H2O.ai, DataRobot). On-platform feature store and Snowflake Marketplace for ML models.
Languages/Interfaces	Python, Scala, SQL, R notebooks; Delta Live Tables; REST APIs. Collaborative notebooks.	ANSI SQL as first-class; also supports Python/Java via Snowpark, Snowflake APIs, and JDBC/ODBC. Now supports Python UDFs and Streams/Tasks for pipelines.
Data Sharing & Collaboration	Delta Sharing for secure data publication; separate workspaces per team. More manual cross-org sharing (storage or custom APIs).	Market-leading data sharing: Secure Data Sharing and Data Marketplace allow live share of data across accounts/orgs without copying (venturebeat.com). Native cross- account collaboration.
Concurrency / Performance	High parallelism on Spark; not as granular for BI concurrency. Autopilot clusters can burst. Good for heavy batch/streaming loads.	Designed for massive concurrency: multiple virtual warehouses can serve many BI/SQL workloads simultaneously without contention. Auto-suspends idle clusters to save cost.
Governance / Security	Delta Lake ACID compliance; RBAC via Unity Catalog; HIPAA-compliant security profile option (docs.databricks.com) (www.databricks.com). Databricks supports major compliance (SOC2, HIPAA) (docs.databricks.com).	Built-in encryption, role-based controls, and enterprise certifications (SOC2, HIPAA). Offers data masking, row-level security, Data Classification tags, etc. Time Travel for compliance audits.
Integration / Ecosystem	Open-source libraries (Delta Lake, MLflow, Spark). Many connectors (Kafka, AWS services, FHIR, etc). Partnerships with genomics/Al tools (Glow, Graphster) (www.prnewswire.com) (www.databricks.com).	Rich partner network (Streamlit, Snowpark, H2O, Dataiku, etc). Works with Tableau, PowerBI, Python H2O modules, etc. Ecosystem focus on data pipelines and analytics apps.

Table 2: Comparison of Databricks and Snowflake platform features (with representative citations).

Data Types and Schemas

Life sciences applications involve highly diverse data:

- Genomic and omics data: Raw sequences (FASTQ), alignments (BAM), variant calls (VCF), gene expression matrices, etc (often petabytes for large projects).
- Imaging and unstructured data: Radiology/PET scans (DICOM), pathology slide images, medical text notes, PDF reports.
- Clinical and phenotypic data: Electronic Health Records (EHR) tables, LIMS outputs, patient metadata, clinical trial data (often semi-structured or relational).
- Research publications and knowledge graphs: Textual and relational data (e.g. publications, ontologies like MeSH, chemical/pathway networks).
- **IoT and sensor data:** Wearable device signals, lab instrument logs (time-series).

Databricks natively handles this multi-modal variety. Its schema-on-read approach and support for unstructured binary (e.g. images) makes it ideal for ingesting "massive volumes of data" that life sciences labs generate (www.databricks.com). The Delta Lake format allows defining metadata (schema) on top, facilitating SQL queries or ML on the same data lake. For example, a hospital could land EHR tables, genomic files, and streaming IoT data all into the Databricks lakehouse, then perform interactive SQL or data science queries. Databricks partners like Glow (genomics library), Graphster (knowledge graphs) provide additional genomic/biomedical data processing (www.prnewswire.com) (www.databricks.com). The Lakehouse for Healthcare product explicitly advertises open-source genomics acceleration ("Glow" for genomics pipelines) as a feature (www.prnewswire.com).

Snowflake, while historically oriented to structured data, has expanded support. Its VARIANT data type supports storing JSON and semi-structured data, but very large files (like images or raw FASTQ) are typically handled via external stages or object storage references. Snowflake recently introduced support for items like Snowflake Python UDFs and external functions to process more complex data, but large unstructured binaries are generally better left outside Snowflake. In short, Snowflake excels when the data can be loaded or transformed into its columns. Common use cases: loading processed genomics results tables, clinical trial data sets, biochemical assay results, etc into Snowflake tables for analytics. For raw high-throughput data, organizations often use Databricks or other cloud storage, then aggregate outcomes to Snowflake for reporting. Snowflake's data exchanges and marketplace can share standard biomedical datasets (e.g. OMOP real-world evidence tables) across companies.

Expert Opinions: A healthcare data expert notes that Databricks' unified platform "bridges the gap between clinical expertise and data-driven insights," enabling end-to-end flows from data ingestion to model building (medium.com). In contrast, Snowflake's analytics environment is viewed as "SQL-first." One industry commentator highlights Snowflake's Time Travel (a feature to query past snapshots of data) as a game-changer for compliance and retrospective analysis in healthcare (medium.com).

Compute, Analytics, and AI/ML

Query Processing and BI Workloads

Snowflake's core is a SQL data warehouse. It is often the first choice for *business intelligence and reporting* in life sciences – e.g. analyzing sales data, financial metrics, clinical trial KPIs, or aggregating real-world evidence. Its automatic clustering and columnar storage give excellent performance for large-scale queries and dashboards. Snowflake's multi-cluster architecture allows many BI teams to concurrently query the same data without slowdown. Snowflake also has natively integrated tools (e.g. Snowpipe for continuous ingestion, Streams/Tasks for change data capture) that make ELT pipelines straightforward.

Databricks also supports SQL querying via *Databricks SQL* and optimized storage (Photon engine in DB SQL Compute). However, because Databricks was built atop Spark, its greatest strength is in **massive-scale data engineering and ML pipelines**. Databricks is commonly used to transform and prepare data at scale before any Bl. It shines when data sizes or complexity exceed what a pure warehouse can easily handle. In many life sciences deployments, a common pattern is to use Databricks to perform intensive processing (e.g. variant calling, bulk genomic analytics, complex physician decision-tree logic) and then write the results into Snowflake for reporting and business analysis.

Machine Learning and AI

Machine learning (ML) is central in modern life sciences: drug discovery increasingly uses AI, and hospitals use predictive models for patient risk and operational efficiency. Both platforms now offer ML capabilities, but with different emphases.

• Databricks: Databricks provides a *unified analytics platform* for data engineering and ML. It includes MLflow (for experiment tracking and model registry), integrates smoothly with popular ML libraries (TensorFlow, PyTorch, scikit-learn), and allows GPUs in cloud clusters. Data scientists can develop code in notebooks (Python/Scala/R) that access the entire data lake directly. Databricks also supports real-time streaming analytics (Spark Structured Streaming) for life sciences, e.g. monitoring patient vitals or lab processes and reacting in real time (medium.com). Internally, Databricks launched an open LLM ("Dolly") and recently acquired MosaicML (for model training) as part of a push to enable generative AI on the platform (techcrunch.com). In the context of healthcare, Databricks advertises "AI use cases like disease prediction, medical image classification, and biomarker discovery" as being enabled by its Lakehouse (www.prnewswire.com). For example, one early adopter (GE Healthcare) reports using Databricks to unify patient data and apply ML for insights (www.prnewswire.com).

• Snowflake: Snowflake's ML story centers on Snowpark, which allows writing Python or Java code that runs within Snowflake's engine. Snowflake also supports external ML through partner integrations: for instance, Snowflake's marketplace includes packages from H2O.ai (for Al model building) and DataRobot, etc. Additionally, Snowflake recently acquired Streamlit to facilitate building data applications, and integrated NVIDIA's NeMo for LLM workloads (techcrunch.com) (techcrunch.com). This means customers can train or fine-tune models using data in Snowflake, often offloaded to GPU resources via partnerships. Snowflake's approach is to let users do ML without leaving the data cloud environment. In life sciences, Snowflake cites use cases like genomics-based predictive models and healthcare predictive analytics on large data sets (medium.com). Snowflake's built-in features like automatic data clustering and materialized views also assist moderate ML tasks by accelerating feature queries.

Case Study (Snowflake + ML): A Snowflake Builders blog highlights an example healthcare scenario: building a demand forecasting model for prescription drugs using Snowpark ML (medium.com). By running Python ML code inside Snowflake. While this is primarily a retail/analytics case, it shows Snowflake's ML in practice. (Snowflake also touts that clients use their platform for RWE modeling and genomics, although detailed case studies are scarce in public literature.)

Performance and Scalability

Benchmarks of raw performance vary widely by workload. Vendor claims and independent tests are mixed. Anecdotally, Databricks often reports faster ETL or complex transformation jobs due to Spark's parallelism, whereas Snowflake shines on pure SQL aggregations. For example, Databricks claims its SQL engine "up to 5x faster and 4x cheaper" on ETL tasks compared to Snowflake in one medium-post benchmark (medium.com) (though that was on specific gueries). Independent analyses suggest Snowflake can outperform on queries with heavy indexing optimization, whereas Databricks outperforms on broad-scan large data extracts. In life sciences use cases, performance also depends on data layout: e.g. Snowflake's micro-partitions autocluster on commonly filtered fields, which helps when querying patient cohorts by age or condition; Databricks can manually optimize partitioning (or use Delta's Z-order) for genomics data by, say, chromosome and position. Concurrency is another factor: Snowflake typically supports many more simultaneous users without contention, while Databricks concurrency depends on how many clusters or jobs run in parallel.

One notable performance metric comes from Snowflake's customer Pfizer: after migrating to Snowflake, Pfizer reported processing data 4× faster and saving 19,000 annual hours (www.snowflake.com) (www.snowflake.com). (A Pfizer executive stated that Snowflake "unified business units with greater access to insights" while cutting cost of ownership by ~57% (www.snowflake.com) (www.snowflake.com).) On the Databricks side, while we lack such quantified public benchmarks, Databricks partners claim major speedups: for instance, a consultancy advertises "4x faster insight generation across science and business functions" and "70% cloud cost optimization through lakehouse and serverless architecture" when using such

platforms (www.agilisium.com). These figures, though promotional, hint at the efficiency gains organizations aim for with modern data stacks.

Data Sharing and Ecosystem

Life sciences research often requires cross-organization data sharing (e.g. clinical trial consortiums, pharmaceutical partners, or public/private data exchanges). Here, Snowflake's **Secure Data Sharing** is a strong point: any Snowflake customer can share live data objects (tables/views) with another Snowflake account without copying, enabling collaboration on sensitive datasets with end-to-end governance (venturebeat.com). Snowflake actively promotes a collaborative healthcare ecosystem: its HCLS Data Cloud is backed by a "live, connected ecosystem" of partner platforms and data providers (venturebeat.com). For example, Snowflake's HCLS platform already counts Anthem, IQVIA, Novartis, Roche among its users (venturebeat.com), demonstrating industry uptake. Snowflake also has a Public Data Sharing program (e.g. OMOP real-world evidence models, common data sets) that can benefit pharma R&D and epidemiological studies.

Databricks historically has less emphasis on cross-company sharing; sharing often occurs via standard open formats (e.g. publishing tables to an S3 bucket or Delta Sharing, or syncing results to Snowflake). However, Databricks is introducing features like Delta Sharing (an open protocol for sharing Delta Lake tables) which is gaining traction among data communities. Databricks' ecosystem focuses more on open-source and data science: for instance, the Lakehouse can integrate partner libraries (Genomics libraries like Glow (www.prnewswire.com), NLP tools like Spark NLP). Databricks also acquired or partnered with key players (e.g. MLflow came from Databricks, Graphster for knowledge graphs) to support health-specific use cases.

Medical informatics initiatives often explore *knowledge graphs* to link disparate biomedical data. Databricks has collaborated on **Life Sciences Knowledge Graphs** built on its platform (www.databricks.com). For example, one solution collates data using biomedical ontologies (MeSH terms, clinical trial metadata) in Delta Lake, enabling graph analytics with SPARQL (www.databricks.com). This approach highlights Databricks' strength in blending graph and AI on its lakehouse. Snowflake can also store graph data in tables, but native graph analytics are not yet a core feature (though third-party graph solutions can run on Snowflake's compute).

Compliance and Security

Data privacy and regulatory compliance are paramount in life sciences. Both Databricks and Snowflake highlight enterprise-grade security:



- Snowflake: The HCLS Data Cloud explicitly "ensures the security, governance and compliance required to meet industry regulations" (venturebeat.com). Snowflake employs strong encryption, network isolation, and auditing. It is HIPAA eligible and compliant with SOC2, ISO, and other standards. Snowflake's built-in features like Time Travel and Data Masking help with auditability and data governance. For example, Snowflake's "Time Travel" lets users query historical data versions, aiding regulatory audits and retrospective compliance (medium.com).
- Databricks: Databricks supports HIPAA workloads via a compliance security profile which enforces encryption, logging, hardened compute images, and requires a Business Associate Agreement (docs.databricks.com). Databricks is SOC2 Type II certified and supports HIPAA, HITRUST, and FedRAMP environments. Its Unity Catalog provides fine-grained access controls and lineage tracking to meet governance needs. The platform also offers "regulatory-grade MLOps" to ensure reproducibility and auditability of Al models (www.databricks.com). For example, Databricks allows keeping clinical trial data and audit trails in a unified platform, addressing GxP (Good Practice) requirements in pharma.

In practice, life science IT teams must architect for compliance whether using Databricks or Snowflake. Both platforms relieve much of the operational burden: one Databricks doc notes that AWS or Azure manages the infrastructure security layers, while Snowflake similarly abstracts the hardware security. Ultimately, both provide strong compliance foundations, but AWS/Azure responsibility models still require customers to configure properly (e.g. enable Databricks' HIPAA profile (docs.databricks.com) or follow Snowflake's best practices).

Use Cases in Life Sciences

This section explores concrete scenarios and success stories for Databricks and Snowflake in life sciences, across the drug lifecycle, healthcare operations, and R&D.

Drug Discovery and Genomics Research

Massive genomic datasets and Al-driven drug discovery define modern biotech. Databricks specifically targets this: its Lakehouse for Healthcare and Life Sciences (launched Mar 2022) was explicitly built for R&D workflows (www.prnewswire.com). The press release (Databricks-Lifesciences Lakehouse) highlights use cases like disease risk prediction, digital pathology classification, real-world evidence ingestion, and using genomics libraries (Glow) for biomarker discovery (www.prnewswire.com). Early adopters like Regeneron and Thermo Fisher Scientific have leveraged the Lakehouse: for instance, a Thermo Fisher IT Director reports that Databricks' platform "enabled us to eliminate costly data silos, unlock new opportunities to innovate, and become a more data-driven organization." (www.prnewswire.com).

Graph-based Knowledge Discovery: The collaboration on a life-sciences knowledge graph (Wisecube + Databricks) shows how Databricks unifies research literature, clinical trials, and genomic annotations. In this example, raw R&D data is loaded into Delta Lake, semantic

ontologies (MeSH, trial keywords) create a graph, and Spark-based graph analytics (via Graphster) run on it (www.databricks.com). The blog argues this allows researchers to ask complex questions (e.g. which gene variants have been linked to a drug's indication) by querying across connected data, something old data silos could not do. This showcases Databricks' strength in handling unstructured search and advanced analytics on biomedical data.

Genomics on Snowflake: Although no Snowflake-native genomics tool exists, the Genes (Basel) 2024 analysis of Snowflake (pmc.ncbi.nlm.nih.gov) describes using Snowflake for genomic data warehousing. It outlines a framework where variant and phenotype tables are loaded into Snowflake for queries, demonstrating "disease variant analysis and in silico drug discovery" using SQL and Python (pmc.ncbi.nlm.nih.gov). The authors show Snowflake can be "convenient and effective" for analyzing cohorts at scale when properly structured. This suggests Snowflake is viable for large-scale genomics when the data is preprocessed into tables.

In summary: For large-scale sequence analysis and ML on genomics, Databricks often leads due to its compute power and genomics libraries (www.prnewswire.com). Snowflake is a strong backend for storing and querying aggregate results (gene counts, biomarkers) and linking structured biomedical databases. Many organizations use Databricks for heavy lifting (sequence alignment, variant calling, complex ML training) and Snowflake for downstream analytics (SQL cohort queries, filtering by phenotype).

Clinical Analytics and Real-World Evidence

Healthcare providers and payers use data platforms for clinical analytics (EHR data), hospital operations, and real-world evidence (RWE) studies. Both Databricks and Snowflake have offerings in this space.

Snowflake in Clinical Data: Snowflake highlights use cases in *patient data integration* and *hospital analytics*. The Snowflake Builders blog mentions implementing Snowflake for RWE and genomics data, suggesting it's already used by teams for COVID-19 and HCP/claims data (medium.com). For example, Snowflake partner accounts like IQVIA and Komodo Health (patient-level data analytics companies) rely on Snowflake's scalable SQL to handle large claims and medical datasets (venturebeat.com). Snowflake's ability to easily join disparate tables (e.g. linking patient claims, prescriptions, lab results) and share schemas facilitates building patient registries or outcomes research. Its multi-cluster architecure suits concurrent research cohorts.

Databricks in Clinical Data: Databricks is used in healthcare analytics for tasks like predictive risk modeling and hospital operations. The Databricks Life Sciences page touts "patient insights at population scale" (www.databricks.com). For instance, one major healthcare provider built a "Health Lakehouse" on Databricks to unify patient records, imaging, and device data, enabling predictive modeling of readmission risk. (Citiations from blog: "Databricks helps organizations build a complete view into patient health" (www.databricks.com).) Databricks also emphasizes

MLOps and reproducibility, which is critical for models that may be subject to FDA review or clinical validation.

Covid-19 Example: During the pandemic, many institutions needed rapid analytics. Snowflake points to customers like hospitals using their platform to analyze COVID-19 testing and vaccination data at scale (medium.com). Databricks likewise offers case studies (e.g. Regeneron used Databricks for genomic discovery at Biobank scale) and talks about enabling rapid collaboration for emergent healthcare data.

Pharmacovigilance/Supply Chain: Life sciences companies also use data platforms for supply chain and regulatory reporting. For example, maintaining a dynamic gross-to-net drug pricing model or pharmacovigilance pipeline. Snowflake's data cloud can store these transactional and IoT (building sensor) data, enabling near-real-time dashboards. A case in point: one vendor case study described using Snowflake + Tableau for prescription drug market analytics (www.thorogood.com). Databricks might handle heavy simulation or forecasting (e.g. optimizing manufacturing yield via ML on sensor data), though concrete references are scarce.

Enterprise Analytics and Business Use Cases

Aside from R&D/medical uses, life sciences companies need enterprise analytics (sales/marketing, finance, supply chain, HR, etc.):

- Databricks can serve as the central analytics engine for complex ETL and ML-driven forecasts. For instance, Agilisium's marketing suggests Databricks builds "ML pipelines, multi-omics processing, and collaborative R&D analytics" in pharma (www.agilisium.com). Databricks' real-time capabilities might support a live "supply chain control tower" model (predict stock-outs, optimize inventory) via streaming R interface.
- Snowflake often dominates static Bl: e.g. field force performance, sales forecasts, financial reporting. Its simplicity is a plus: as one blog noted for a retail customer, Snowflake's simplicity and performance for BI led that company to choose Snowflake for its analytics (www.devtechie.com). Similarly, in pharma/medtech, Snowflake clusters could run credit system queries or personalizing HCP profiling. Snowflake's separation of compute also allows separate departments (commercial, manufacturing, regulatory) to run heavy queries on the same data without resource contention.

Both platforms are increasingly being combined in industry solutions. For example, a consulting firm describes a framework where Snowflake provides a real-time analytic data layer (for trial data, commercial metrics, claims) while Databricks handles R&D ML pipelines (e.g. ingesting, cleaning, modeling for predictive insights) (www.agilisium.com). This hybrid approach leverages each system's strengths (Databricks for heavy lifting AI, Snowflake for democratized analytics).

Case Studies and Examples

Pfizer - Snowflake: Pfizer migrated large parts of its analytics to Snowflake. According to Snowflake, Pfizer "processes data 4x faster" with Snowpark and achieved a 57% reduction in total cost of ownership (www.snowflake.com) (www.snowflake.com). Snowflake claims Pfizer saved ~19,000 labor hours annually by speeding up data tasks. In this transition, Pfizer unified multiple disparate business units onto Snowflake, enabling cross-unit data sharing and faster reporting.

Regeneron - Databricks: Regeneron, a genomics-focused pharma, uses Databricks for largescale genome analysis. In a webinar, Regeneron's team described running bulk variant analysis pipelines and leveraging Databricks for interactive data exploration on biobank-scale data (www.databricks.com). (Databricks even produced a detailed report: "How Regeneron Accelerates Genomic Discovery at Biobank Scale".) While not Snowflake-specific, Regeneron's use underlines Databricks' strengths in genomics R&D at scale.

Thermo Fisher Scientific - Databricks: In the PR release (www.prnewswire.com), Thermo Fisher (scientific instrument maker) credited Databricks with helping them "eliminate costly data silos" across research data. Their Sr. IT Director called Databricks "a modern platform for data and AI" that unlocks innovation. This likely refers to Thermo's analytics environment integrating lab data across units (e.g. combining genomic assay results, inventory data, QA logs) to accelerate product development.

Anthem, Novartis, IQVIA - Snowflake: Snowflake's venture into healthcare lists major adopters. For example, Anthem (a large insurer) is reported using Snowflake's platform to centralize claims and clinical data for analytics (venturebeat.com). Global pharma companies like Novartis and Roche are also listed as Snowflake customers, indicating use for both research and enterprise analytics. (Specifics aren't given, but one can infer they use it for global data consolidation and regulated analytics.)

AMN Healthcare (Migrating off Databricks): In one LinkedIn report, AMN Healthcare moved from Databricks to Snowflake, achieving 93% lower data lake costs, a pipeline success rate of 99.9%, and 75% reduction in warehouse runtime. (However, LinkedIn posts are not peerreviewed sources; we note it only as an industry anecdote that some firms do re-platform flexibility.)

These examples illustrate that success with either platform depends on aligning to the workload. Pfizer's dramatic improvement on Snowflake suggests large structured data burdens can be streamlined there. Regeneron's genomic tasks needed Databricks' parallelism. Realworld pilots often use both: e.g. ingest data in Databricks, analyze in Snowflake, or share curated data back into Snowflake for Bl.

Data Analysis and Evidence-Based Arguments



To evaluate Databricks vs Snowflake from an evidence-based standpoint, we consider metrics like cost, performance, time-to-insight, and usability, as reported in case studies or benchmarks.

- TCO and Efficiency: Vendor claims aside, independent analysis points to Snowflake's lower overhead for pure SQL workloads. Its pay-per-use model (compute paused when idle) can save costs for sporadic workloads. Databricks, using Spark clusters, may be more costly if clusters run 24/7 for streaming or modeling. However, Databricks encourages multi-tenancy and auto-scaling to optimize. In life sciences with continuous pipelines (e.g. nightly ETL from EHRs), Snowflake's auto-suspension can be especially cost-effective.
- Skillsets: Data teams in pharma often are either IT/SQL-centric or data science-centric. Surveys suggest business analysts prefer Snowflake's familiar SQL environment, whereas data scientists gravitate to Databricks' notebook interfaces and Python/R support. This splits the user base: Snowflake may be easier to adopt for BI/devs already skilled in SQL, while Databricks can be more productive for predictive modeling.
- Platform Convergence: Notably, both companies are extending into each other's territory. Databricks has bolstered its SQL query performance (Photon engine), and Snowflake has expanded beyond warehousing into data engineering (Snowpark, Python UDFs) and ML support (techcrunch.com). Also, each is embracing AI: by mid-2023, Snowflake acquired Neeva/Streamlit and partnered with NVIDIA to enable in-platform LLM training and inference (techcrunch.com) (techcrunch.com). Meanwhile, Databricks open-sourced "Dolly" (an LLM) and acquired MosaicML for integrated model building (techcrunch.com). These moves suggest future convergence: both platforms are becoming AI-first data clouds. For life sciences, this means that either platform can serve as a hub for advanced analytics and generative AI on biomedical data.
- Ecosystem Lock-In: Another factor is ecosystem compatibility. Databricks being closely tied to Spark and open-source may feel more "portable" (projects can be moved to other Spark-compatible platforms if needed). Snowflake, while enabling data exports, is a proprietary service; migrating off Snowflake requires effort. Life sciences companies must weigh this when investing.

Case Studies and Real-World Examples

We explicitly incorporate additional case narratives:

• BioPharma R&D (Databricks & Snowflake hybrid): One large pharmaceutical firm built an intelligent R&D platform using both tools. They used Snowflake as a secured data warehouse for curated experiment and clinical datasets (e.g. compound screening results, patient registry data), enabling easy reporting and regulatory audit. Databricks complemented this by handling unstructured research data (raw omics, image processing) and hosting notebooks for crossdisciplinary teams. This dual-platform approach leveraged Snowflake's governance for final datasets and Databricks' power for heavy analytics. (Company names are rarely disclosed in public sources due to confidentiality, but such multi-platform strategies are known in consulting circles.)

- Retail & Supply Chain in Life Sciences (Snowflake): A drug distributor used Snowflake plus Tableau to create a prescription analytics dashboard (www.thorogood.com). By loading sales, claims, and marketing data into Snowflake, and linking it to Tableau's visualization, they achieved near-realtime insight into market dynamics. This suggests Snowflake's strength in supply-chain and commercial analytics for pharma.
- Thermo Fisher & Medical Device (Databricks): A global life sciences equipment manufacturer ingrained Databricks in its digital strategy. A senior executive from the company stated Databricks Lakehouse "helps accelerate precision medicine and fundamentally change care by predicting disease" (www.prnewswire.com). Internally, they likely use Databricks to analyze sensor data from instruments, streamline genomics workflows, and deliver Al-enhanced software to customers. The exact ROI metrics aren't public, but the strong endorsement indicates how an instrumentation company saw databricks as core to its innovation roadmap.
- Consultant Claim (Agilisium): An industry consultancy advertises that deploying modern data lakes and warehouses yields "4× faster insight generation" and "70% cloud cost optimization" (www.agilisium.com). While marketing, this underscores the expected orders-of-magnitude improvement in data analytics agility life sciences teams seek. However, such metrics will vary; real results depend on initial baselines.

Discussion: Strategic and Future Implications

The choice between Databricks and Snowflake is not strictly a winner-takes-all. In life sciences, many organizations employ both platforms for complementary purposes. Key considerations:

- Breadth vs Depth of Workloads: Databricks is often better for experimental R&D analytics (genomics, high-throughput screening, image analysis, real-time monitoring) due to its flexibility and compute power (www.prnewswire.com) (www.databricks.com). Snowflake tends to suit standardized analytics and cross-organization data sharing (commercial analytics, multi-site clinical studies, enterprise BI) (venturebeat.com) (medium.com). Companies with heavy multi-modal R&D data may tilt toward Databricks; those prioritizing enterprise data consolidation may lean on Snowflake.
- Talent and Process: The platforms support different workflows. Databricks encourages notebookdriven, agile development often resembling a data science environment. Snowflake focuses on managed warehouse processes, which might integrate with traditional ETL tools or SQL-based pipelines. Organizations may need to train staff on both sets of skills.
- Governance and Compliance: Both can meet regulatory requirements, but Snowflake's managed nature means fewer customer-managed components. Life sciences organizations regulated by FDA or EMA might prefer Snowflake's simplicity for validated environments, although Databricks has a path with its compliance security profile (docs.databricks.com). Auditors tend to scrutinize histrionic pipelines; having a clear chain-of-custody (more straightforward in a single SQL warehouse) can be an advantage.

- Vendor Ecosystem: Future roadmaps will influence decisions. Databricks' partnership with OpenAl (making Al agents available inside Databricks via "Agent Bricks") and Snowflake's with Nvidia (GPUs for in-cloud AI) (techcrunch.com) (techcrunch.com) indicate that each is becoming a platform for enterprise AI. For life sciences, this means easier integration of advanced AI tools (e.g. generative models for drug design and text mining) directly where data resides.
- Data Collaboration: The trend in life sciences is toward networked data (cooperative research, data consortia). Snowflake's vision of a Data Cloud, where multiple stakeholders share data within a governed environment, aligns with this. For instance, pharma companies are exploring sharing anonymized clinical data via covalent platforms - Snowflake's marketplace could facilitate this. Databricks, with its open paradigm, might participate via open protocols (e.g. Delta Sharing) and align with initiatives like the Global Alliance for Genomics and Health (GA4GH).

Future Directions: The evolving landscape suggests both platforms will continue to add lifesciences-specific features. Possible directions include:

- · Genomic Data Warehousing: Building specialized genomic-optimized tables or query engines (post-SQL on sequences) - Snowflake could introduce genomic data types, while Databricks might integrate more bioinformatics tools.
- Federated Learning: For privacy-sensitive health data, federated analytics (train across silos without centralizing data) are emerging. Both platforms could support federated queries or model training (Snowflake's Secure Data Sharing, Databricks federated clusters, etc).
- · AI/ML Integration: As noted, both companies see AI as core. Expect deeper integration of LLMs and bio-Al: e.g. pre-built medical LLM on top of Snowflake data, or Databricks offering more MLOps automation tailored to FDA pipelines.
- Hybrid/On-Prem Support: Some research data remains on-prem (e.g. vulnerable patient data in hospitals). Snowflake now offers a hosted option in Azure Stack, and Databricks runs on AWS Outposts or Azure (and even on-prem Spark). Solutions that bridge on-prem and cloud will matter to cautious organizations.
- Standardization: As more life science software is deployed in the cloud, standards (FHIR for health records, OMOP for observational data) become important. Weather both Databricks or Snowflake natively support these models out-of-the-box will influence their adoption. Currently, both support loading/saving these schemas, but deeper tool integration (e.g. FHIR endpoints, built-in OMOP pipelines) could be areas of future investment.

Conclusion

The comparison between Databricks and Snowflake in life sciences reveals no universal winner rather, each platform excels in different facets of the domain. Databricks' lakehouse architecture, built for data science and AI, is ideal for cutting-edge R&D (genomics, imaging, predictive modelling) where handling petabyte-scale multi-modal data is crucial (www.prnewswire.com) (www.databricks.com). Snowflake's cloud data platform, with its

powerful SQL engine and sharing capabilities, is well-suited for enterprise analytics, clinical data warehousing, and collaboration across organizations (venturebeat.com) (www.snowflake.com). In practice, top life sciences firms often deploy both: for instance, using Databricks to lay the foundation for large-scale machine learning and Snowflake to democratize the refined data for business use.

Our survey of literature, case studies, and technical resources shows that life sciences organizations should align platform choice with workflow requirements. Key considerations include data types (unstructured genomics vs. structured clinical), workloads (batch ML pipelines vs. live dashboards), compliance needs, and organizational expertise. Both Databricks and Snowflake are rapidly evolving: their recent pushes into generative AI, knowledge engineering, and cross-cloud data ecosystems suggest that future capabilities will increasingly overlap. For example, Snowflake's HCLS Data Cloud and Databricks' Healthcare Lakehouse were launched within a week of each other, underscoring how both companies aim to be the "onestop data platform" for healthcare and pharma (venturebeat.com) (www.prnewswire.com).

Recommendations: Life sciences companies embarking on a data platform strategy should conduct thorough workload analyses. If bioinformatics and AI/ML innovation are the priority, investing in Databricks and its data science tooling may yield the fastest time-to-innovation. If regulatory compliance, multi-source integration, and BI are paramount, Snowflake's governed warehouse may reduce time-to-production. In many cases, a hybrid approach will prevail. Wherever possible, organizations should leverage both platforms' strengths in a cohesive architecture.

In conclusion, Databricks and Snowflake are transforming life sciences IT. By intelligently applying each platform's strengths to the domain's unique data challenges—such as integrating patient-genomic profiles or scaling clinical trials analytics—life science enterprises can accelerate discovery, improve patient outcomes, and maintain a competitive edge in an Al-driven future (medium.com) (www.prnewswire.com).

IntuitionLabs - Industry Leadership & Services

North America's #1 Al Software Development Firm for Pharmaceutical & Biotech: IntuitionLabs leads the US market in custom Al software development and pharma implementations with proven results across public biotech and pharmaceutical companies.

Elite Client Portfolio: Trusted by NASDAQ-listed pharmaceutical companies including Scilex Holding Company (SCLX) and leading CROs across North America.

Regulatory Excellence: Only US AI consultancy with comprehensive FDA, EMA, and 21 CFR Part 11 compliance expertise for pharmaceutical drug development and commercialization.

Founder Excellence: Led by Adrien Laurent, San Francisco Bay Area-based AI expert with 20+ years in software development, multiple successful exits, and patent holder. Recognized as one of the top AI experts in the USA.

Custom Al Software Development: Build tailored pharmaceutical Al applications, custom CRMs, chatbots, and ERP systems with advanced analytics and regulatory compliance capabilities.

Private Al Infrastructure: Secure air-gapped Al deployments, on-premise LLM hosting, and private cloud Al infrastructure for pharmaceutical companies requiring data isolation and compliance.

Document Processing Systems: Advanced PDF parsing, unstructured to structured data conversion, automated document analysis, and intelligent data extraction from clinical and regulatory documents.

Custom CRM Development: Build tailored pharmaceutical CRM solutions, Veeva integrations, and custom field force applications with advanced analytics and reporting capabilities.

Al Chatbot Development: Create intelligent medical information chatbots, GenAl sales assistants, and automated customer service solutions for pharma companies.

Custom ERP Development: Design and develop pharmaceutical-specific ERP systems, inventory management solutions, and regulatory compliance platforms.

Big Data & Analytics: Large-scale data processing, predictive modeling, clinical trial analytics, and real-time pharmaceutical market intelligence systems.

Dashboard & Visualization: Interactive business intelligence dashboards, real-time KPI monitoring, and custom data visualization solutions for pharmaceutical insights.

Al Consulting & Training: Comprehensive Al strategy development, team training programs, and implementation guidance for pharmaceutical organizations adopting Al technologies.

Contact founder Adrien Laurent and team at https://intuitionlabs.ai/contact for a consultation.

DISCLAIMER

The information contained in this document is provided for educational and informational purposes only. We make no representations or warranties of any kind, express or implied, about the completeness, accuracy, reliability, suitability, or availability of the information contained herein.

Any reliance you place on such information is strictly at your own risk. In no event will IntuitionLabs.ai or its representatives be liable for any loss or damage including without limitation, indirect or consequential loss or damage, or any loss or damage whatsoever arising from the use of information presented in this document.

This document may contain content generated with the assistance of artificial intelligence technologies. Al-generated content may contain errors, omissions, or inaccuracies. Readers are advised to independently verify any critical information before acting upon it.

All product names, logos, brands, trademarks, and registered trademarks mentioned in this document are the property of their respective owners. All company, product, and service names used in this document are for identification purposes only. Use of these names, logos, trademarks, and brands does not imply endorsement by the respective trademark holders.

IntuitionLabs.ai is North America's leading AI software development firm specializing exclusively in pharmaceutical and biotech companies. As the premier US-based AI software development company for drug development and commercialization, we deliver cutting-edge custom AI applications, private LLM infrastructure, document processing systems, custom CRM/ERP development, and regulatory compliance software. Founded in 2023 by Adrien Laurent, a top AI expert and multiple-exit founder with 20 years of software development experience and patent holder, based in the San Francisco Bay Area.

This document does not constitute professional or legal advice. For specific guidance related to your business needs, please consult with appropriate qualified professionals.

© 2025 IntuitionLabs.ai. All rights reserved.