

DNA Gene Mapping Project

**An Integrated Genomics, Disease, and Variant
Intelligence Platform**

Prepared by

Sharique Mohammad

February 2026

Copyright © 2026 Sharique Mohammad

All rights reserved.

This document, including its structure, data models, processing logic, feature engineering strategies, and analytical workflows, is the intellectual property of the author. No part of this document may be reproduced, distributed, or transmitted in any form or by any means without prior written permission, except for personal, academic, or evaluative review purposes with proper attribution.

The software designs, data schemas, and analytical methodologies described herein are provided solely for educational and research purposes and are not intended for clinical, diagnostic, or commercial use.

Abstract

This project presents the design and implementation of a comprehensive DNA Gene Mapping platform that integrates genomic variants, gene annotations, disease knowledge, protein information, and clinical relevance into a unified analytical system. The platform is built using a layered data architecture that includes raw ingestion, structured enrichment, feature engineering, and advanced analytics readiness.

Public biomedical data sources including ClinVar, MedGen, OMIM, RefSeq, UniProt Swiss-Prot, and related reference datasets are systematically ingested, normalized, and linked to construct high-quality gene, variant, and disease representations. The project emphasizes data integrity, traceability, and reproducibility through explicit schema design, feature lineage tracking, and modular processing pipelines.

Beyond foundational data modeling, the system is engineered to support advanced downstream use cases such as disease association discovery, polygenic risk analysis, gene prioritization, and readiness for deep learning applications including convolutional neural networks, graph-based models, and transformer architectures.

The entire platform is implemented with a local-first execution strategy, prioritizing accessibility, cost efficiency, and computational feasibility in constrained hardware environments. This work demonstrates how genomics analytics and machine learning readiness can be achieved without reliance on large-scale cloud infrastructure.

Regulatory and Usage Disclaimer

This project is intended solely for research and educational purposes. It is not designed, validated, or approved for clinical diagnosis, medical decision-making, patient treatment, or regulatory submission.

All biological, genomic, and disease-related data used in this project originate from publicly available sources. Although care has been taken to preserve the accuracy, provenance, and consistency of the data, the author makes no guarantees regarding the completeness, correctness, or real-world applicability of the information presented.

No clinical interpretations, therapeutic recommendations, or patient-level decisions should be derived from this system. Users are solely responsible for ensuring compliance with applicable laws, ethical standards, and institutional guidelines when applying, extending, or interpreting this work.

Contents

Abstract	ii
Regulatory and Usage Disclaimer	iii
1 Executive Overview	1
1.1 Project Vision	1
1.2 Target Users (Research, Clinical, Pharma)	2
1.3 Why This Project Matters	3
2 Architectural Motivation and Design Principles	5
2.1 Motivation: Fragmentation of Genomic Data	5
2.2 Limitations of Existing Public Databases	5
2.3 Engineering Versus Research Gaps	6
2.4 Design Principles	6
3 High-Level System Architecture	8
3.1 End-to-End Architecture Overview (Conceptual)	8
3.2 Layered Data Architecture Philosophy	9
3.3 Role of Databricks Community Edition	10
3.4 Why Silver and Gold Layers Are the Analytical Focus	10
4 Data Ingestion, Persistence, and Schema Governance	11
4.1 Execution Environment and Local-First Constraints	11
4.2 Spark and PySpark as the Analytical Backbone	11
4.3 Table Persistence Strategy	12
4.4 Schema Discipline and Governance	12
4.5 Preparation for Analytical and Machine Learning Workflows	12
5 Data Source Acquisition and Governance	13
5.1 ClinVar Ingestion	13
5.2 MedGen Ingestion	13
5.3 OMIM Integration	14
5.4 RefSeq References	14
5.5 UniProt Swiss-Prot	14
5.6 Genetic Testing Registry (GTR)	14
5.7 Data Licensing Considerations	15
6 Gene Modeling and Enrichment	16

6.1	Gene Identifier Normalization	16
6.2	Alias Resolution Strategy	18
6.3	Functional Annotation Extraction	18
6.4	Druggability and Pharmacogene Tagging	18
6.5	Genomic Coordinate Handling	18
7	Variant Processing Pipeline	19
7.1	Variant Normalization	19
7.2	Allele Mapping Logic	20
7.3	Clinical Significance Parsing	21
7.4	Transcript and Protein Consequence Parsing	21
7.5	Quality Scoring Logic	21
8	Disease Modeling and Integration	22
8.1	MedGen Concepts	22
8.2	Disease Hierarchy	23
8.3	OMIM and MONDO Alignment	24
8.4	Disease Enrichment Logic	24
8.5	Gene–Disease Link Construction	24
9	Protein and Transcript Layer	25
9.1	RefSeq Transcripts	25
9.2	Protein Mapping Strategy	26
9.3	UniProt Annotations	27
9.4	Protein–Gene–Variant Linkage	27
10	Feature Engineering Framework	28
10.1	Silver-to-Gold Transformation Rules	28
10.2	Functional Impact Feature Construction	29
10.3	Population Frequency Integration and Leakage Prevention	30
10.4	Feature Grouping and Confidence Scoring	30
11	Classical Machine Learning	31
11.1	Disease-Level Feature Construction	31
11.2	Gene Prioritization Feature Framework	32
11.3	Gene–Disease Association Scoring	32
11.4	Evaluation and Validation Signals	33
11.5	Explainability and Transparency	33
12	Deep Learning and Graph Modeling	34
12.1	Why Deep Learning Is Optional but Valuable	34
12.2	Model-Agnostic Feature Preparation	34
12.3	Convolutional Neural Networks for Sequence-Based Modeling	35
12.4	Graph Neural Networks for Biological Networks	35
12.5	Transformer-Based Contextual Modeling	35

12.6	Bias Awareness and Mitigation Strategies	35
12.7	Execution Constraints and Architectural Planning	36
12.8	What Is Safe to Implement Now vs Later	36
13	Power BI Analytics	37
13.1	Dashboard Architecture	37
13.2	Gene Analytics Dashboards	37
13.3	Variant Analytics Dashboards	38
13.4	Disease Risk and Complexity Dashboards	38
13.5	Local Execution and Usage Constraints	38
14	Advanced Analytics and Visualization Layer	39
14.1	Analytical Objectives	39
14.2	Visualization Consumption Strategy	39
14.3	Dashboard Design Themes	40
14.4	Data Refresh and Governance	40
15	Machine Learning Model Implementation and Evaluation	41
15.1	Execution Environment	41
15.2	Initial Modeling Use Cases	41
15.3	Model Selection Philosophy	42
15.4	Evaluation Strategy	42
16	System Architecture and Application Layer Design	43
16.1	Backend Architecture	43
16.2	Frontend Architecture	43
16.3	Project Structure and Domain Separation	44
16.4	Security and Access Control	44
17	Backend API Layer	45
17.1	FastAPI Architecture	45
17.2	API Endpoint Design	45
17.3	Security Considerations	46
17.4	Model Serving Strategy	46
18	Web Application Frontend	47
18.1	React Architecture	47
18.2	Component Structure	47
18.3	API Integration	48
18.4	User Experience Principles for Genomics	48
19	System Integration and Deployment	49
19.1	End-to-End System Flow	49
19.2	Execution and Orchestration Strategy	50
19.3	Monitoring and Observability	50
19.4	Scalability Considerations	50

19.5	Future Portability	51
20	ML Feature Inventory (High-Level)	52
20.1	Gene-Level Features	52
20.2	Variant-Level Features	53
20.3	Disease-Level Features	53
20.4	Cross-Entity and Derived Features	54
20.5	Feature Governance Principles	54
21	Feature Lineage Documentation	55
21.1	Source-to-Feature Traceability	55
21.2	Lineage Across Data Layers	55
21.3	Auditing and Reproducibility	56
21.4	Explainability by Design	56
22	Conclusion	57
	Acknowledgement of Data Sources	58
	References	59
	Appendices	60
A	Data Layer Schemas	61
A.1	Reference Layer Table Schemas	62
A.2	Silver Layer Table Schemas	62
A.3	Gold Layer Table Schemas	63
B	Feature Dictionary (Core Biological Entities)	64
B.1	Gene Features (Silver Layer)	64
B.2	Variant Features (Silver Layer)	64
B.3	Disease Features (Silver Layer)	65
B.4	Gene–Disease Association Features (Silver Layer)	66
B.5	Disease Machine Learning Features (Gold Layer)	66
C	Feature Dictionary (Derived & Analytical Features)	68
C.1	Disease Association Features	68
C.2	Polygenic Risk Features	69
C.3	Gene Prioritization Features	69
D	Feature Lineage Documentation	70
D.1	Feature Lineage Principles	70
D.2	Example Feature Lineage	70
D.3	Lineage Enforcement Strategy	71
D.4	Lineage Scope and Limitations	72
E	System Architecture Reference (Logical)	73
E.1	Logical Architecture Layers	73
E.2	Technology Stack Summary	74
E.3	Architectural Scope and Boundaries	75

F	Governance and Compliance Reference	76
F.1	Regulatory Disclaimer	76
F.2	Ethical Safeguards	76
F.3	Data Licensing and Attribution	77
F.4	Governance Enforcement Strategy	77
F.5	Scope Boundaries and Intended Use	77
G	System Architecture (Local Execution)	79
G.1	End-to-End System Architecture (Local)	79
G.2	Data Layer Architecture (Logical, Not Cloud-Based)	79
G.3	Feature Lineage and Dependency Tracking	80
G.4	Gene–Variant–Disease Relationship Model	80
H	Machine Learning Architecture (Local Research Mode)	82
H.1	Architectural Scope and Design Principles	82
H.2	End-to-End Machine Learning Flow	83
H.3	Model Categories Supported	83
H.4	Local Execution Constraints	84
H.5	Feature Persistence and Governance	84
I	Analytics and Visualization Layer (Local)	85
I.1	Analytical Consumption Strategy	85
I.2	Power BI Desktop Architecture	86
I.3	Dashboard Design Principles	86
I.4	Data Refresh and Governance	86
I.5	Local-Only Execution Constraints	87
J	Deep Learning Architecture (Conceptual)	88
J.1	Architectural Scope and Design Principles	88
J.2	Conceptual Deep Learning Flow	89
J.3	Neural Model Families and Conceptual Roles	89
J.4	Feature-Oriented Integration Philosophy	90
J.5	Local Execution Constraints	90
K	Deep Learning Execution Architecture	91
K.1	Execution Overview	91
K.2	CNN Execution Pipeline	92
K.3	GCN Execution Pipeline	93
K.4	Transformer Execution Pipeline	94
K.5	Feature Integration Rules	95
K.6	Execution Constraints	96
L	Model-to-Feature Mapping Inventory	97
L.1	Purpose and Mapping Principles	97
L.2	CNN-Derived Feature Mapping	98

L.3	GCN-Derived Feature Mapping	98
L.4	Transformer-Derived Feature Mapping	99
L.5	Feature Governance Rules	100
M	Training, Evaluation, and Validation Strategy	101
M.1	Training Strategy	101
M.2	Evaluation Metrics	102
M.3	Validation Approach	102
M.4	Risk and Limitations	103
M.5	Final Integration Rule	103
N	Local Execution Constraints and Design Decisions	104
N.1	Hardware Constraints	104
N.2	Software Constraints	105
N.3	Data Volume Management Strategy	105
N.4	Compute-Aware Feature Engineering	105
N.5	Deferred Infrastructure Components	106
N.6	Research-First Architecture Justification	106
O	Project Folder Structure and Codebase Organization	107
O.1	Top-Level Repository Structure	107
O.2	Machine Learning and Deep Learning Workspace	108
O.3	Backend Application Structure	108
O.4	Frontend Application Structure	109
O.5	Deployment and Environment Separation	110

List of Figures

1.1	High-level system architecture illustrating the Bronze–Silver–Gold data layers and downstream analytical readiness.	2
3.1	Conceptual end-to-end architecture of the DNA Gene Mapping Project, showing the layered progression from raw data ingestion (Bronze) through normalized entities (Silver) and feature-engineered analytical outputs (Gold).	8
6.1	Gene-centric enrichment workflow illustrating identifier normalization, alias resolution, functional annotation, and druggability tagging leading to the curated Silver-layer gene table.	17
7.1	Variant processing and enrichment workflow showing coordinate normalization, allele mapping, clinical significance parsing, consequence extraction, and quality scoring leading to the curated Silver-layer variant table.	20
8.1	Disease ontology alignment and integration flow showing MedGen as the primary disease entity with OMIM and MONDO integrated as cross-references to support normalized variant–disease and gene–disease relationships.	23
9.1	Protein–gene–variant integration model showing genes as the central biological anchor, RefSeq transcripts as reference mappings, and UniProt proteins as curated functional entities linked to genomic variation.	26
10.1	Silver-to-Gold feature engineering workflow showing the integration of evolutionary conservation, functional impact, population frequency, and clinical evidence to produce stable Gold-layer variant features suitable for downstream analytics and machine learning.	29
11.1	Classical machine learning feature framework showing the aggregation of Gold-layer variant, gene, and disease features into structured feature sets, followed by scoring and prioritization to produce interpretable, machine-learning-ready inputs.	32

A.1	Logical relationship between reference, silver, and gold data layers. The reference layer provides authoritative normalization and controlled vocabularies, the silver layer models normalized biological entities, and the gold layer contains analytical and machine-learning-ready feature tables derived from silver-layer data.	61
H.1	Logical machine learning architecture in local research mode. Curated Silver and Gold layer features feed classical, graph-based, and deep learning models. All models operate in a CPU-only environment and persist outputs as reusable analytical features rather than final predictions.	83
J.1	Conceptual deep learning architecture showing how curated Silver and Gold tables feed CNN, GCN, and Transformer models. Neural outputs are extracted as structured features and persisted into Gold-layer feature tables under local, CPU-only execution constraints.	89
K.1	Deep learning execution flow showing curated Silver and Gold tables feeding CNN, GCN, and Transformer architectures, with learned representations persisted as interpretable Gold-layer features under local execution constraints.	92
K.2	CNN execution pipeline illustrating sequence encoding, convolutional feature extraction, pooling, and dense layers used to derive interpretable variant- and protein-level feature representations.	93
K.3	GCN execution flow showing curated graph construction, neighborhood aggregation, node embedding generation, and disease–gene risk signal derivation.	94
K.4	Transformer execution pipeline showing tokenization, attention-based contextual embedding, and feature extraction for downstream analytical use.	95
O.1	Top-level monorepo structure showing separation between data, machine learning, backend services, frontend applications, and deployment artifacts.	108
O.2	Machine learning and deep learning workspace structure separating feature engineering, classical ML, graph-based models, deep learning, training, and evaluation.	109
O.3	Backend service structure illustrating application modules, routing, schema enforcement, service layers, and deployment artifacts.	110
O.4	Frontend application structure showing modular components, services, hooks, type definitions, and build configuration.	111
O.5	Deployment directory structure containing container orchestration, reverse proxy configuration, and environment templates.	112

List of Tables

1.1	Target user groups and supported analytical capabilities	3
A.1	Reference layer core tables and primary identifiers	62
A.2	Silver layer core biological entity tables	63
A.3	Gold layer analytical and machine learning feature tables	63
B.1	Core gene features (silver layer)	65
B.2	Core variant features (silver layer)	65
B.3	Core disease features (silver layer)	66
B.4	Gene–disease association features (silver layer)	66
B.5	Disease-level machine learning features (gold layer)	67
C.1	Derived disease association features	68
C.2	Derived polygenic risk features	69
C.3	Derived gene prioritization features	69
L.1	CNN model outputs mapped to analytical features	98
L.2	GCN model outputs mapped to analytical features	99
L.3	Transformer model outputs mapped to analytical features	99

Chapter 1

Executive Overview

1.1 Project Vision

The DNA Gene Mapping Project is a comprehensive biomedical data science and machine learning system designed to integrate genomic variants, genes, proteins, and disease knowledge into a unified analytical platform. The core vision of the project is to move beyond fragmented bioinformatics scripts and academic notebooks and instead build an end-to-end system that reflects how real-world clinical, pharmaceutical, and translational research platforms are engineered.

Modern genomic data are inherently complex and distributed across multiple authoritative sources, such as ClinVar, MedGen, OMIM, RefSeq, UniProt, gnomAD, and dbVar. Individually, these datasets are valuable, but in isolation they do not provide a complete or analytically usable representation of biological reality. This project is designed to unify these disparate sources into a coherent, queryable, and machine-learning-ready system.

The platform emphasizes layered data modeling, including a raw ingestion layer (Bronze), structured normalization and enrichment (Silver), and feature-engineered analytical outputs (Gold). Strict schema enforcement, feature lineage tracking, and clear separation of concerns between ingestion, enrichment, analytics, and modeling are central design principles. Rather than optimizing for speed of development or visual presentation, the system prioritizes correctness, traceability, explainability, and long-term extensibility.

The layered data foundation underlying this platform, including the separation of reference, silver, and gold data representations, is formally defined in Appendix A.

Figure 1.1 provides a high-level view of the system architecture and the layered data model used throughout the project.

High-Level System Architecture

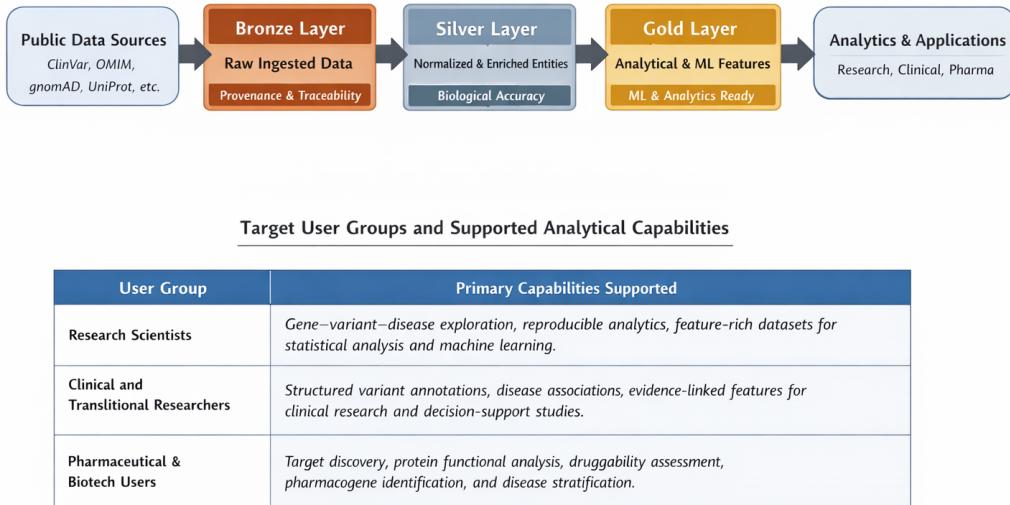


Figure 1.1: High-level system architecture illustrating the Bronze–Silver–Gold data layers and downstream analytical readiness.

1.2 Target Users (Research, Clinical, Pharma)

The project is designed to serve multiple overlapping user profiles within the biomedical and life sciences ecosystem.

For researchers, the platform enables systematic exploration of gene–variant–disease relationships, reproducible analytics, and feature-rich datasets suitable for statistical analysis and machine learning experimentation. Researchers benefit from normalized identifiers, resolved aliases, and curated disease mappings that significantly reduce the preprocessing effort.

For clinical and translational researchers, the system provides structured representations of clinically relevant variants, disease associations, and evidence-linked annotations. Although not intended for direct diagnostic use, the platform mirrors the internal data preparation layers commonly found in hospital research environments and clinical decision support research systems.

For pharmaceutical and biotech users, the project demonstrates how genomic and proteomic data can be transformed into analytically useful assets for target discovery, druggability analysis, pharmacogene identification, and disease stratification. The inclusion of protein annotations, functional features, and disease complexity metrics aligns with real-world drug discovery and biomarker research workflows.

Table 1.1 outlines the primary user groups and the corresponding analytical capabilities supported by the platform.

User Group	Primary Capabilities Supported
Research Scientists	Systematic exploration of gene–variant–disease relationships, reproducible analytical workflows, and feature-rich datasets suitable for statistical analysis and machine learning experimentation.
Clinical and Translational Researchers	Structured representations of clinically relevant variants, disease associations, and evidence-linked annotations to support clinical research, validation studies, and decision-support system development (non-diagnostic).
Pharmaceutical and Biotech Users	Transformation of genomic and proteomic data into analytically useful assets for target discovery, protein functional analysis, druggability assessment, pharmacogene identification, and disease stratification.

Table 1.1: Target user groups and supported analytical capabilities

1.3 Why This Project Matters

This project addresses three fundamental challenges that limit the practical use of genomics data.

First, biological and clinical data are highly fragmented across heterogeneous sources, each with its own identifiers, schemas, and semantics. Without careful integration, this fragmentation prevents holistic analysis.

Second, raw genomic data are not directly usable for analytics or machine learning. Extensive normalization, enrichment, and semantic alignment are required before meaningful insights can be derived. Many public examples stop at ingestion and do not demonstrate this critical transformation layer.

Third, many academic or tutorial-based projects do not show how genomic data can be operationalized for downstream analytics, disease modeling, and predictive intelligence in a system-oriented setting.

This project addresses all three challenges simultaneously. It demonstrates how authoritative public datasets can be transformed into a structured, feature-rich analytical system that supports advanced use cases while remaining feasible under real-world constraints.

The system is intentionally developed under limited hardware, storage, and budget conditions. These constraints are treated as design drivers rather than limitations, rein-

forcing the project’s focus on efficient engineering, disciplined architecture, and practical feasibility. The result is a realistic blueprint for building robust genomics analytics platforms without reliance on enterprise-scale infrastructure.

Regulatory scope, ethical boundaries, and intended usage constraints for this project are formally documented in Appendix F.

Chapter 2

Architectural Motivation and Design Principles

2.1 Motivation: Fragmentation of Genomic Data

Genomic and biomedical data have grown exponentially over the past two decades. Advances in sequencing technologies have enabled the identification of millions of genetic variants across populations, yet interpreting those variants in a biologically and clinically meaningful way remains a major challenge. A single human genome may contain millions of variants, while only a small fraction are associated with known functional or clinical relevance.

Biomedical knowledge is distributed across multiple authoritative public resources, each curated with a specific scope and purpose. For example, ClinVar focuses on clinical interpretations of variants, OMIM catalogs Mendelian disease associations, MedGen provides disease concept normalization, UniProt curates protein structure and function, gnomAD captures population allele frequencies, and dbVar stores information on structural variants.

While each dataset is valuable in isolation, none provides a complete or analytically usable representation of biological reality on its own. Meaningful insight emerges only when these heterogeneous resources are integrated, reconciled, and analyzed together within a unified analytical system.

2.2 Limitations of Existing Public Databases

Despite their scientific value, existing public genomic databases present several limitations when used directly for analytics or machine learning.

First, identifier fragmentation is pervasive. Genes, variants, and diseases are represented using multiple identifier systems, naming conventions, and aliases across data

sources. The same biological entity may appear under different identifiers or labels, making naive joins unreliable and error-prone without careful normalization.

Second, schema inconsistency poses a significant engineering challenge. Each data source defines its own schema structure, field semantics, and data types. Raw ingestion without structured transformation leads to brittle pipelines, ambiguous columns, and analytically unusable tables.

Third, genomic data are characterized by extreme scale and sparsity. Variant-level datasets are large, while most individual variants carry limited annotation. This imbalance complicates downstream analysis and increases the risk of misleading conclusions if not handled systematically.

Finally, public datasets are not designed to be directly consumed by analytical or machine learning models. Raw genomic and clinical annotations must be transformed into carefully engineered features that preserve biological meaning while remaining computationally tractable.

The structured biological entities introduced to address these limitations are formally defined in Appendix B.

2.3 Engineering Versus Research Gaps

The core challenge addressed by this project is not sequence alignment, variant calling, or primary annotation. Instead, it focuses on the post-annotation problem: transforming annotated genomic data into a structured, queryable, and model-ready analytical system.

Many academic or tutorial-based genomics projects emphasize exploratory analysis or isolated modeling tasks, often relying on ad-hoc scripts and implicit assumptions. While valuable for research experimentation, such approaches rarely demonstrate how genomic data can be operationalized within a disciplined engineering framework that supports reproducibility, traceability, and long-term extensibility.

This project explicitly targets the engineering gap between raw biomedical data and downstream analytics. It aims to bridge research-oriented exploration and system-oriented design by enforcing consistent schemas, preserving lineage from raw ingestion to engineered features, and enabling complex analytical and machine learning workflows without sacrificing transparency or auditability.

The analytical and derived feature constructs introduced to bridge these gaps are documented in Appendix C.

2.4 Design Principles

The architecture of the DNA Gene Mapping Project is guided by a set of explicit design principles that inform all downstream implementation decisions.

Correctness and biological validity take precedence over convenience or development speed. All transformations are designed to preserve semantic meaning and minimize ambiguity.

Traceability is enforced across the entire data lifecycle. Data lineage is maintained from raw ingested datasets through normalized entities and feature-engineered outputs, enabling reproducibility and auditability.

Separation of concerns is strictly applied. Raw ingestion, normalization and enrichment, analytical feature generation, and modeling are treated as distinct stages with clearly defined responsibilities.

Schema discipline is a central requirement. Explicit, versioned schemas are used to prevent silent drift, reduce coupling between components, and support long-term maintainability.

Finally, the system is designed under realistic execution constraints. Limited hardware, storage, and budget considerations are treated as architectural drivers rather than limitations, encouraging efficient engineering practices and practical feasibility.

Chapter 3

High-Level System Architecture

3.1 End-to-End Architecture Overview (Conceptual)

Figure 3.1 presents a conceptual, end-to-end view of the system architecture, illustrating the progression from raw data ingestion through normalized biological entities and feature-engineered analytical outputs.

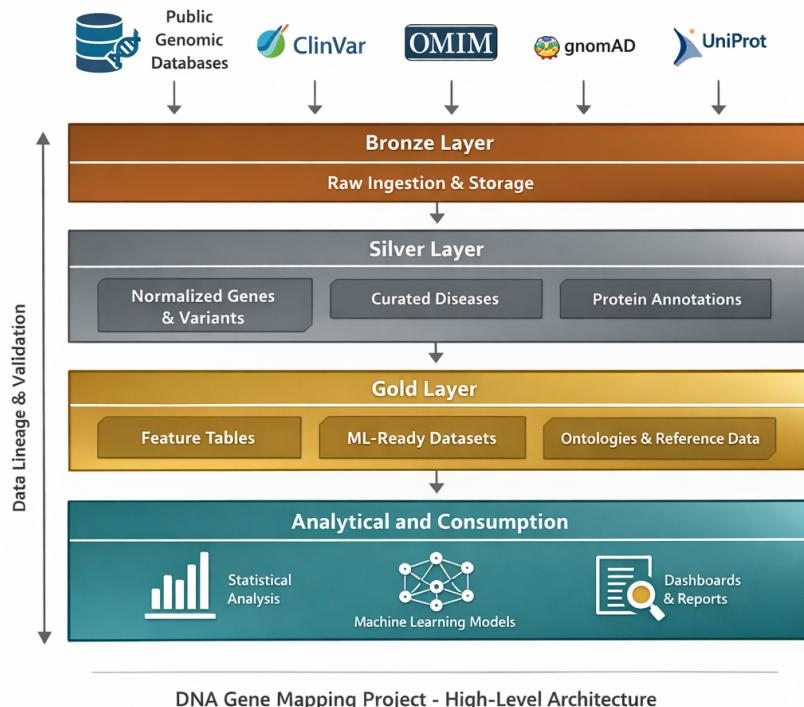


Figure 3.1: Conceptual end-to-end architecture of the DNA Gene Mapping Project, showing the layered progression from raw data ingestion (Bronze) through normalized entities (Silver) and feature-engineered analytical outputs (Gold).

At a high level, the DNA Gene Mapping Project is structured as a layered analyt-

ical system that transforms heterogeneous public biomedical datasets into structured, biologically consistent, and analytically usable representations. The architecture emphasizes a clear separation of responsibilities between data ingestion, normalization, feature engineering, and downstream consumption.

Rather than treating the system as a single monolithic pipeline, the architecture is designed as a sequence of well-defined layers. Each layer produces stable outputs that serve as explicit inputs to subsequent stages, enabling traceability, reproducibility, and controlled system evolution.

A formal logical breakdown of architectural layers is provided in Appendix E.

3.2 Layered Data Architecture Philosophy

The system follows a layered data architecture inspired by modern analytical platforms used in both industry and research environments. This design minimizes coupling between components, reduces technical debt, and ensures that changes in upstream data sources do not unpredictably propagate into downstream analytics.

Conceptually, the architecture is organized into the following layers:

- **Bronze layer (raw ingestion)** — minimally processed datasets ingested directly from public biomedical sources. This layer preserves raw provenance and supports reprocessing but is not intended for direct analytical use.
- **Silver layer (normalized entities)** — cleaned, normalized, and biologically enriched representations of core entities such as genes, variants, diseases, and proteins. This layer enforces consistent identifiers, resolved aliases, and semantic alignment across sources.
- **Gold layer (analytical features)** — feature-engineered tables derived from Silver entities, designed explicitly for analytics, statistical modeling, and machine learning consumption.
- **Reference layer** — stable lookup, ontology, and normalization tables shared across layers to ensure consistency and controlled vocabulary usage.

This layered separation ensures that analytical logic operates only on curated, semantically stable data and that feature definitions remain reproducible across system iterations.

The local execution realization of this layered architecture is described in Appendix E.

3.3 Role of Databricks Community Edition

Databricks Community Edition serves as the primary analytical execution environment for this project. It is used strictly as a distributed data processing and analytics engine, not as a managed cloud platform or architectural dependency.

Within the system architecture, Databricks is responsible for:

- Executing distributed transformations using Spark and PySpark
- Enforcing schema-aware table creation and validation
- Supporting iterative analytical development under local execution constraints
- Enabling scalable feature engineering without reliance on external infrastructure

The architectural concepts described in this chapter are intentionally independent of Databricks itself. Databricks is an implementation choice that operationalizes these concepts while respecting the project's local-first execution strategy and hardware limitations.

3.4 Why Silver and Gold Layers Are the Analytical Focus

Although the architecture includes a conceptual Bronze layer for raw data ingestion, analytical workflows in this project deliberately operate only on the Silver, Gold, and Reference layers.

The Silver layer provides biologically consistent, normalized entities that are suitable for reliable querying, enrichment, and relationship modeling. The Gold layer builds upon Silver outputs to generate stable, model-consumable features with explicit lineage and controlled semantics.

This design choice ensures that:

- Raw ingestion changes do not propagate directly into analytics
- Feature definitions remain stable and reproducible
- Machine learning models consume curated representations rather than raw data
- Analytical results remain explainable and auditable

By intentionally restricting analytical access to Silver and Gold layers, the system prioritizes correctness, interpretability, and long-term maintainability over direct exposure of raw datasets.

Chapter 4

Data Ingestion, Persistence, and Schema Governance

4.1 Execution Environment and Local-First Constraints

All data processing and transformation workflows in this project are executed using Databricks Community Edition under a strictly local-first execution model. This environment provides distributed processing capabilities through Apache Spark while remaining compatible with limited hardware, storage, and budget constraints.

The choice of Databricks Community Edition is driven by its ability to support scalable analytical transformations, schema-aware table management, and iterative development without reliance on commercial cloud infrastructure. Importantly, Databricks functions as an execution environment rather than an architectural dependency, implementing the layered design described in earlier chapters.

4.2 Spark and PySpark as the Analytical Backbone

Apache Spark, accessed through PySpark, serves as the primary analytical engine for data transformation and enrichment. PySpark enables expressive transformation logic while preserving performance characteristics required for large-scale joins, aggregations, and feature computation.

Transformation logic is implemented using declarative DataFrame operations wherever possible to ensure readability, debuggability, and deterministic execution. This approach supports reproducibility and aligns with the project's emphasis on transparent data processing.

4.3 Table Persistence Strategy

Intermediate and final datasets are persisted as managed tables within the Databricks environment. This includes all Silver, Gold, and Reference layer tables. Persisted tables act as stable contracts between pipeline stages, preventing repeated recomputation and enabling incremental development.

Raw ingested datasets corresponding to the Bronze layer may be stored transiently or selectively persisted based on size and reuse considerations. Regardless of persistence strategy, all Bronze-layer data are treated as immutable inputs to downstream transformations.

4.4 Schema Discipline and Governance

Explicit schema definition is a foundational design requirement of the system. All persisted tables adhere to predefined schemas that specify column names, data types, primary identifiers, and semantic intent.

Schema enforcement serves multiple purposes:

- Preventing silent schema drift across pipeline iterations
- Enabling reliable downstream analytics and feature engineering
- Supporting traceability and feature lineage documentation
- Ensuring compatibility with analytical and machine learning workflows

Schema definitions are treated as authoritative contracts rather than implementation details. Any schema evolution is deliberate, versioned, and documented to preserve system integrity.

4.5 Preparation for Analytical and Machine Learning Workflows

Although advanced machine learning and deep learning models are introduced in later chapters, data ingestion and persistence decisions are made with downstream analytical requirements in mind.

Silver and Gold tables are structured to support classical statistical analysis, graph-based learning, and neural modeling without requiring reprocessing of raw data. Feature engineering logic prioritizes biological interpretability, controlled dimensionality, and computational feasibility under constrained execution environments.

Chapter 5

Data Source Acquisition and Governance

This chapter corresponds to Phase 1 of the DNA Gene Mapping Project and focuses on the acquisition of authoritative biomedical data sources and the establishment of governance principles. The objective of this phase is to ensure long-term sustainability, traceability, and regulatory awareness across all downstream components of the pipeline.

5.1 ClinVar Ingestion

ClinVar serves as the primary source of clinically interpreted genetic variants within the system. The ingestion process focuses on extracting variant identifiers, allele identifiers, genomic coordinates, clinical significance assertions, review status, submission counts, and associated phenotypes.

Special handling is implemented to support variants with multiple submissions and conflicting interpretations. Original ClinVar values are preserved to maintain provenance, while normalized representations are created to enable standardized downstream analysis. Temporal fields, including the last evaluated date, are retained to support clinical relevance scoring and recency-aware analytics.

5.2 MedGen Ingestion

MedGen is ingested as the authoritative disease ontology and identifier backbone for the project. The MedGen concepts dataset is parsed to extract MedGen identifiers, preferred disease names, source provenance, and semantic classifications.

Only preferred disease terms are retained for analytical modeling to ensure consistency in disease labeling across downstream tables. Synonym and alternative naming information is preserved within reference structures to support enrichment, search, and

cross-database alignment workflows.

5.3 OMIM Integration

OMIM integration is performed using curated mapping resources that link genes to OMIM entries and corresponding MedGen concepts. OMIM identifiers are retained as a complementary clinical reference due to their importance in Mendelian disease modeling and inheritance-based interpretation.

OMIM data is not treated as a standalone disease ontology. Instead, it is aligned with MedGen identifiers to avoid duplication, semantic drift, and inconsistent disease representations across the system.

5.4 RefSeq References

RefSeq reference datasets are ingested to provide authoritative gene and transcript definitions aligned to the GRCh38 genome assembly. RefSeq serves as the foundation for transcript modeling, genomic coordinate normalization, and gene-to-protein relationships.

Extracted attributes include gene identifiers, transcript identifiers, chromosomal coordinates, strand orientation, and reference source metadata. These normalized references support consistent mapping across variant, protein, and disease modeling workflows.

5.5 UniProt Swiss-Prot

UniProt Swiss-Prot is selected as the curated protein knowledge source for the project. Ingestion is restricted to reviewed human protein records, explicitly excluding unreviewed TrEMBL entries to maintain annotation quality and reliability.

Extracted attributes include UniProt accessions, protein names, associated gene symbols, organism taxonomy, functional descriptions, and curated annotations relevant to clinical and pharmaceutical interpretation.

5.6 Genetic Testing Registry (GTR)

The Genetic Testing Registry is ingested to establish links between genes and real-world clinical genetic tests. Extracted fields include test identifiers, test names, associated genes, and disease contexts.

GTR integration provides a direct connection between genomic data and applied clinical diagnostics, reinforcing the translational orientation of the system without introducing patient-level or proprietary information.

5.7 Data Licensing Considerations

All data sources incorporated into the project are evaluated for licensing terms, usage restrictions, and attribution requirements prior to ingestion. Only publicly available datasets with permissive usage conditions are included.

No proprietary, restricted, or patient-identifiable datasets are used. This ensures that the project remains suitable for academic research, educational use, and professional portfolio demonstration without legal or ethical risk.

Formal governance rules, licensing principles, and ethical safeguards are documented in Appendix F.

Chapter 6

Gene Modeling and Enrichment

This chapter corresponds to Phase 2 of the DNA Gene Mapping Project and focuses on transforming raw gene metadata into a richly annotated, normalized, and analytically robust gene representation. The objective of this phase is to establish a stable gene-centric foundation that supports downstream variant analysis, disease modeling, and feature engineering.

6.1 Gene Identifier Normalization

Gene identifiers from multiple authoritative sources are normalized into a unified gene-centric model. Entrez Gene ID is selected as the primary numeric identifier, while HGNC, Ensembl, UniProt, RefSeq, and Alliance Genome identifiers are retained as cross-references.

Figure 6.1 illustrates the gene-centric enrichment workflow, showing how raw gene references from multiple authoritative sources are transformed into a unified Silver-layer gene representation.

Gene-Centric Enrichment Process

DNA Gene Mapping Project

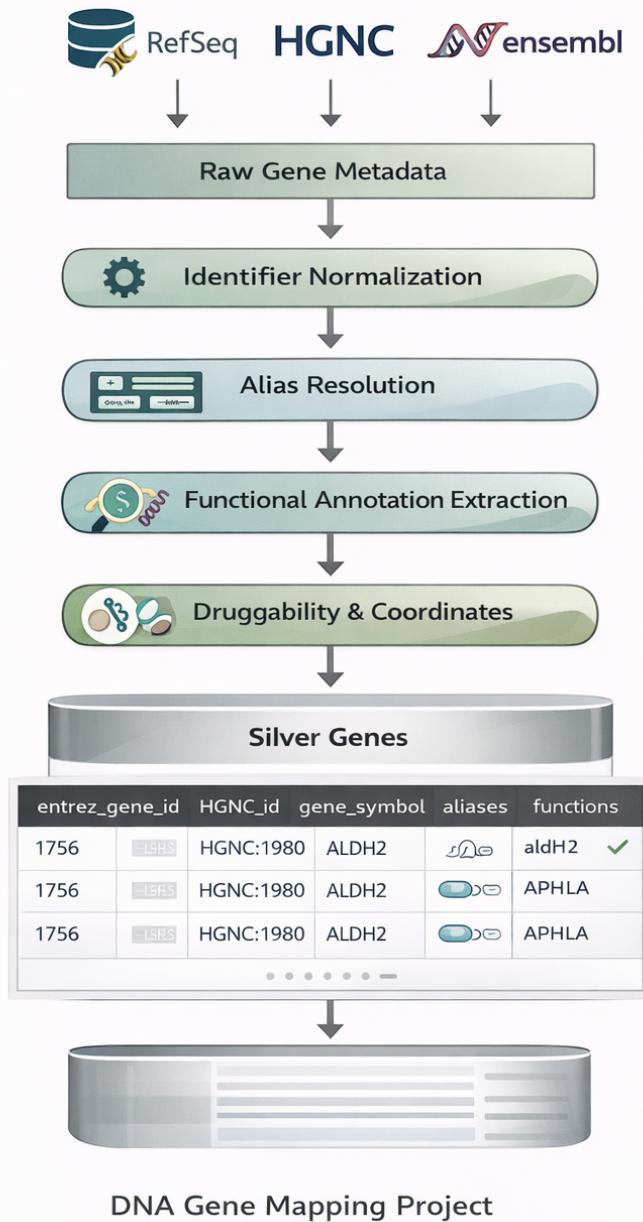


Figure 6.1: Gene-centric enrichment workflow illustrating identifier normalization, alias resolution, functional annotation, and druggability tagging leading to the curated Silver-layer gene table.

This normalization strategy ensures consistent joins across datasets and eliminates ambiguity caused by synonym-based or free-text matching. By enforcing a single primary identifier, the system maintains referential integrity across Silver and Gold layer tables.

The reference and silver-layer table structures supporting gene normalization are defined in Appendix A.

6.2 Alias Resolution Strategy

Gene aliases are extracted, tokenized, and expanded into structured alias representations. A controlled alias resolution strategy is implemented to prevent alias collisions and reduce false-positive joins across datasets.

Alias counts are computed to quantify naming complexity on a per-gene basis. These metrics support downstream confidence scoring and help identify genes with high ambiguity or inconsistent nomenclature across sources.

6.3 Functional Annotation Extraction

Gene-level functional descriptions are parsed to extract biologically relevant annotations, including molecular roles, biological processes, and functional categories. Rule-based natural language heuristics are applied to identify attributes such as enzyme activity, receptor function, transcriptional regulation, and binding properties.

This approach enables functional categorization without reliance on external machine learning models, ensuring interpretability and reproducibility under constrained execution environments.

Core gene feature definitions are enumerated in Appendix B.

6.4 Druggability and Pharmacogene Tagging

Genes are evaluated for druggability and pharmacogenomic relevance using curated keyword heuristics and controlled database cross-references. Boolean flags and categorical indicators are derived to support pharmaceutical and precision medicine analyses.

Representative attributes include indicators for pharmacogene status, known drug targets, cytochrome P450 involvement, and drug metabolism roles. These features enable downstream prioritization of clinically and therapeutically relevant genes.

6.5 Genomic Coordinate Handling

Gene genomic coordinates are normalized to the GRCh38 reference assembly. Derived attributes include gene length, chromosomal location, cytogenetic band assignments, chromosomal arm positions, and proximity to telomeric or centromeric regions.

These spatial attributes support structural genomic analyses and enable consistent integration with variant-level coordinate data in subsequent pipeline stages.

Chapter 7

Variant Processing Pipeline

This chapter corresponds to Phase 3 of the DNA Gene Mapping Project and focuses on building a scalable, interpretable, and schema-consistent variant processing pipeline. The objective of this phase is to transform heterogeneous variant annotations into normalized, analyzable representations suitable for clinical research and downstream feature engineering.

7.1 Variant Normalization

Genetic variants are normalized using consistent genomic coordinate representation, standardized chromosome formatting, and explicit separation of reference and alternate alleles. Normalization logic ensures compatibility with the GRCh38 reference assembly and enables reliable joins with gene reference data and downstream analytical tables.

Figure 7.1 presents the end-to-end variant processing workflow, illustrating how raw variant records are transformed into curated Silver-layer variant entities through normalization, enrichment, and quality scoring.

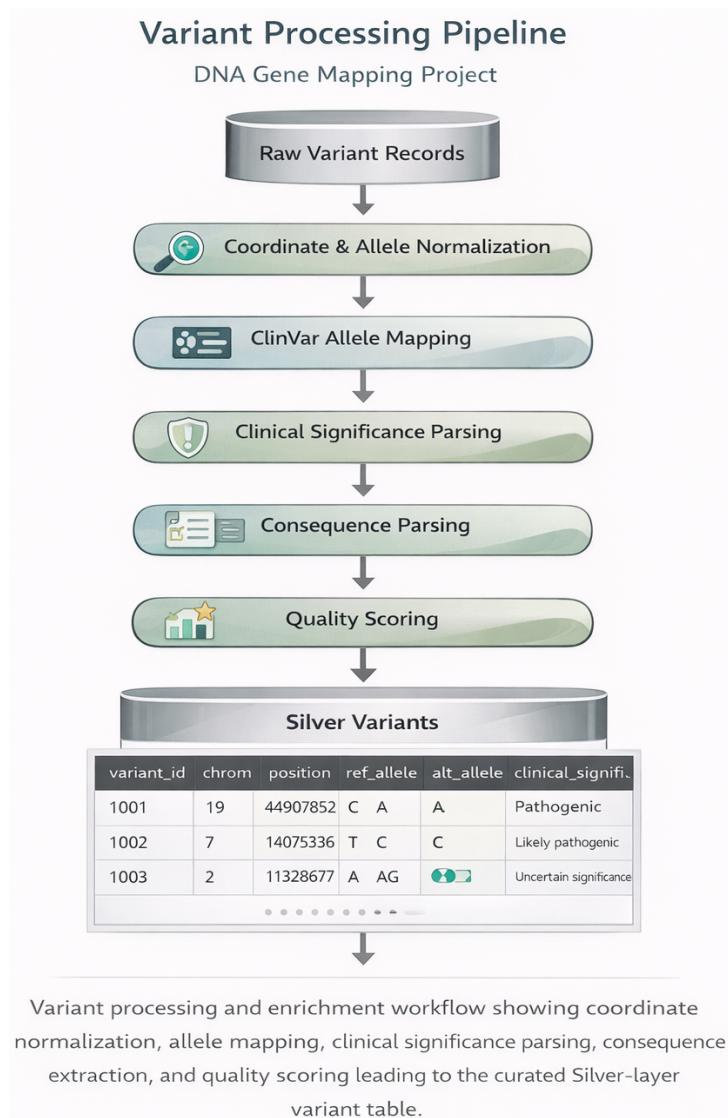


Figure 7.1: Variant processing and enrichment workflow showing coordinate normalization, allele mapping, clinical significance parsing, consequence extraction, and quality scoring leading to the curated Silver-layer variant table.

HGVS-style notations are parsed where available to support interoperability with external annotation services, while original source values are preserved to maintain provenance and traceability.

Variant-level table schemas and data layer responsibilities are summarized in Appendix A.

7.2 Allele Mapping Logic

Allele identifiers are mapped using explicit ClinVar variation and allele relationships. Missing, ambiguous, or conflicting allele mappings are explicitly flagged rather than inferred.

This conservative mapping strategy preserves data integrity, prevents silent enrichment errors, and ensures that downstream analyses can distinguish unresolved mappings from biologically absent information.

7.3 Clinical Significance Parsing

Clinical significance annotations are standardized into controlled categorical values, including pathogenic, likely pathogenic, benign, likely benign, and variants of uncertain significance.

Conflicting interpretations are retained rather than collapsed. Review status, number of independent submitters, and assertion consistency are extracted as structured indicators to support downstream confidence scoring and analytical prioritization.

Variant-level feature definitions and controlled vocabularies are specified in Appendix B.

7.4 Transcript and Protein Consequence Parsing

Variant consequence fields describing transcript-level and protein-level effects are parsed to extract interpretable, schema-aligned attributes. Extracted features include positional indicators, frameshift flags, nonsense and missense classifications, splice-site disruption indicators, and protein impact categories.

No transcript, cDNA, or sequence entities are modeled directly. Instead, consequence strings are decomposed into structured features that enable efficient filtering, aggregation, and machine learning use without introducing sequence-level complexity.

7.5 Quality Scoring Logic

Variant-level quality scores are computed using a combination of review status, number of submitters, and evaluation recency. These scores provide a quantitative confidence signal that complements categorical clinical significance labels.

Quality scores are designed to support downstream analytics and machine learning workflows while remaining transparent, explainable, and reproducible.

Derived analytical features used for downstream disease association and prioritization are defined in Appendix C.

The feature lineage framework ensuring traceability from raw variant annotations to analytical confidence scores is documented in Appendix D.

Chapter 8

Disease Modeling and Integration

This chapter corresponds to Phase 4 of the DNA Gene Mapping Project and focuses on establishing a consistent, ontology-aligned disease modeling layer grounded in real-world clinical semantics. The objective of this phase is to normalize disease representations and enable reliable integration across genes, variants, and downstream analytical features.

8.1 MedGen Concepts

Disease modeling in this project is centered on the need to reconcile multiple biomedical ontologies into a single, analytically consistent representation. Different sources describe diseases using overlapping but non-identical identifiers and hierarchical structures, which can lead to ambiguity if not handled explicitly.

To address this, MedGen is selected as the primary disease entity within the system. Other disease ontologies are aligned with MedGen rather than modeled independently, ensuring semantic consistency and preventing duplication across downstream tables.

Figure 8.1 illustrates this disease modeling and ontology alignment strategy, highlighting MedGen as the central disease representation with OMIM and MONDO integrated as controlled cross-references that support normalized variant–disease and gene–disease relationships.

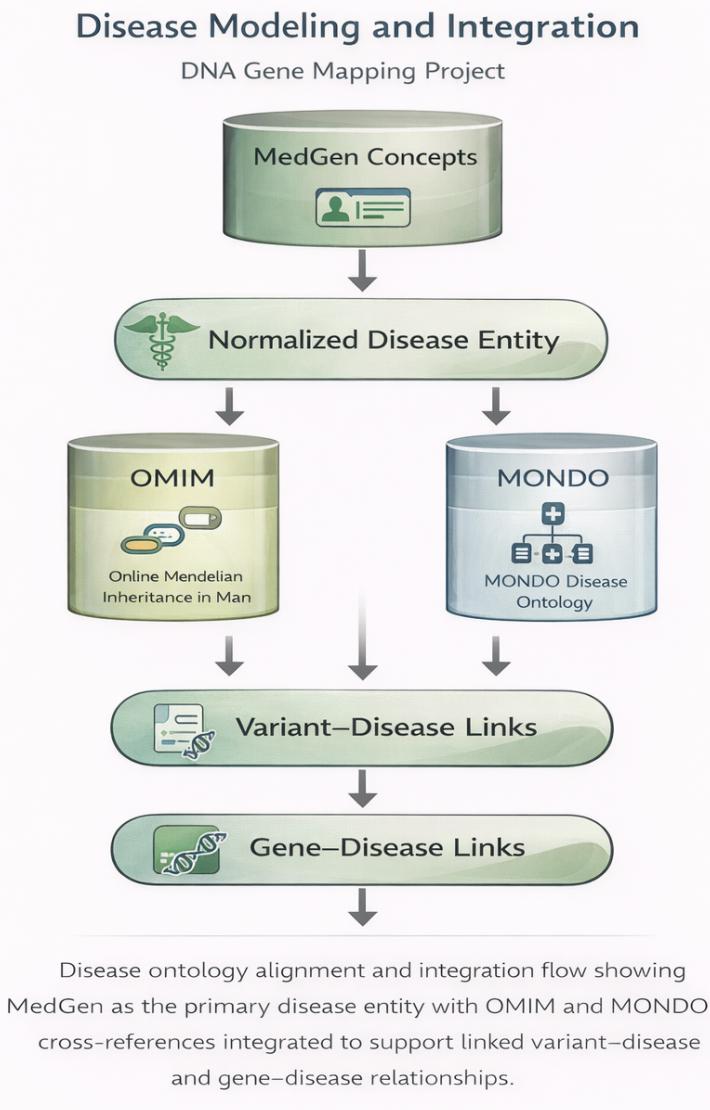


Figure 8.1: Disease ontology alignment and integration flow showing MedGen as the primary disease entity with OMIM and MONDO integrated as cross-references to support normalized variant–disease and gene–disease relationships.

Source provenance is retained for each disease entry, allowing transparent traceability to originating databases and supporting controlled enrichment from complementary resources.

Disease-related table schemas and normalization layers are documented in Appendix A.

8.2 Disease Hierarchy

Hierarchical relationships between disease concepts are preserved to support parent–child disease analysis, phenotype generalization, and disease grouping. These relationships enable analytical queries at multiple levels of disease granularity without duplicating

entities or flattening ontology structure.

Disease hierarchies are treated as reference relationships rather than inferred classifications, ensuring that analytical logic remains aligned with curated biomedical knowledge.

8.3 OMIM and MONDO Alignment

OMIM and MONDO identifiers are aligned to MedGen concepts where authoritative mappings are available. These identifiers are retained as cross-references rather than independent disease entities.

This alignment strategy enables cross-ontology interoperability while preventing duplication, semantic drift, and inconsistent disease representations across the system.

8.4 Disease Enrichment Logic

Variant-level disease labels originating from upstream sources are enriched using validated reference mappings. Generic, ambiguous, or source-specific disease labels are replaced with standardized MedGen concepts whenever possible.

Original disease annotations are preserved for provenance, while normalized disease identifiers are used for analytics and feature engineering to ensure consistency and interpretability.

8.5 Gene–Disease Link Construction

Gene–disease associations are constructed using validated MedGen and OMIM references, resulting in explicit many-to-many relationship mappings. These links support disease burden analysis, polygenic modeling, and downstream network-based analytical workflows.

Gene–disease relationships are treated as first-class analytical entities, enabling systematic exploration of disease complexity without embedding assumptions about causality or clinical actionability.

Disease and gene–disease association feature definitions are provided in Appendix B.

Chapter 9

Protein and Transcript Layer

This chapter corresponds to Phase 5 of the DNA Gene Mapping Project and completes the core biological modeling stack by integrating transcript references and curated protein annotations. The objective of this phase is to enable consistent linkage between genes, variants, transcripts, and proteins while maintaining analytical simplicity and schema discipline.

9.1 RefSeq Transcripts

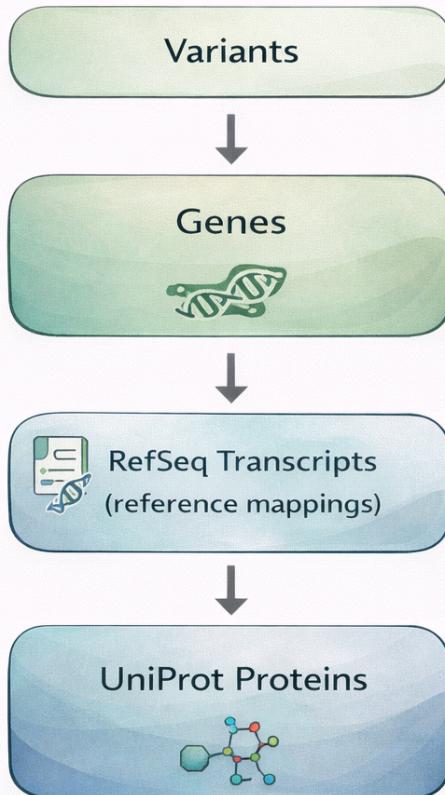
To complete the biological modeling stack, it is necessary to explicitly define how variants, genes, transcripts, and proteins are linked without introducing sequence-level complexity.

RefSeq transcript records are incorporated as structured references to support accurate gene and protein mapping. Extracted attributes include transcript identifiers, associated gene identifiers, genomic coordinates, strand orientation, and reference source metadata.

Figure 9.1 illustrates the integration strategy used in this project, highlighting genes as the central anchor with transcripts serving as reference mappings and proteins providing curated functional context.

Protein–Gene–Variant Integration

DNA Gene Mapping Project



Protein–gene–variant integration model showing genes as the central biological anchor, *RefSeq transcripts* as reference mappings, and *UniProt proteins* as curated functional entities linked to genomic variation.

Figure 9.1: Protein–gene–variant integration model showing genes as the central biological anchor, RefSeq transcripts as reference mappings, and UniProt proteins as curated functional entities linked to genomic variation.

Transcripts are treated as reference-level entities rather than analytical objects. This design avoids transcript-level explosion while preserving the necessary linkage information required for variant consequence interpretation and protein mapping.

Protein and transcript table schemas are defined in Appendix A.

9.2 Protein Mapping Strategy

Protein accessions are mapped to genes using controlled cross-references derived from RefSeq and UniProt. This strategy ensures consistent and reproducible protein–gene

associations across the dataset.

Protein mappings are validated against curated reference sources to prevent ambiguous or indirect associations, enabling stable downstream integration with variant and disease data.

9.3 UniProt Annotations

UniProt Swiss-Prot annotations are integrated to provide high-quality, manually curated protein metadata. Only reviewed human protein entries are included to maintain annotation reliability and clinical relevance.

Integrated attributes include protein names, functional descriptions, biological roles, domain information, and curated comments relevant to disease mechanisms and pharmaceutical research.

9.4 Protein–Gene–Variant Linkage

Final linkage is established between variants, proteins, and genes using normalized identifiers and validated reference mappings. This linkage completes the core biological graph of the system and enables multi-level analytical queries spanning genomic variation, gene function, and protein biology.

These relationships form a critical foundation for downstream disease modeling, feature engineering, and machine learning workflows without introducing unnecessary modeling complexity or sequence-level dependencies.

Protein and transcript-related feature structures are documented in Appendix B.

Chapter 10

Feature Engineering Framework

This chapter corresponds to Phase 6 of the DNA Gene Mapping Project and marks the transition from deterministic data engineering into analytical enrichment. The focus of this phase is the construction of interpretable, schema-governed features derived from curated biological signals such as evolutionary conservation, population frequency, functional impact, and clinical evidence.

10.1 Silver-to-Gold Transformation Rules

Silver-to-Gold feature engineering involves the systematic integration of multiple biological and clinical signals into stable, machine-learning-ready representations. Figure 10.1 provides a high-level view of this transformation, illustrating how curated Silver-layer variant data is enriched with conservation, functional impact, population frequency, and clinical evidence to produce Gold-layer features.

Silver-to-Gold transformations are governed by explicit, reproducible rules designed to preserve biological meaning while producing analytically stable features. Evolutionary conservation scores are integrated to quantify how intolerant genomic positions are to variation.

Rather than relying on locally hosted static databases, conservation data is retrieved through aggregated annotation services. Integrated scores include PhyloP, PhastCons, and GERP, each capturing distinct aspects of evolutionary constraint across vertebrate genomes.

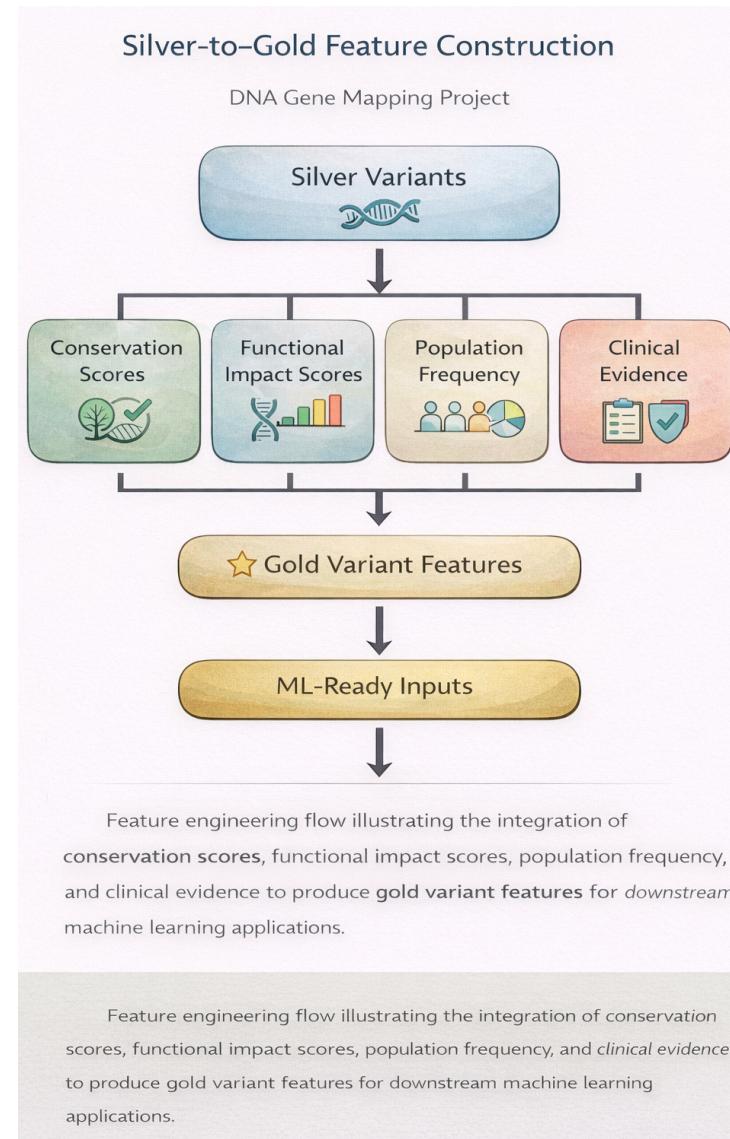


Figure 10.1: Silver-to-Gold feature engineering workflow showing the integration of evolutionary conservation, functional impact, population frequency, and clinical evidence to produce stable Gold-layer variant features suitable for downstream analytics and machine learning.

All conservation scores are normalized and stored at the variant level. Boolean conservation indicators are derived to flag highly conserved, constrained, and evolutionarily significant variants for downstream analysis.

Derived analytical feature groups and transformation outputs are defined in Appendix C.

10.2 Functional Impact Feature Construction

Functional impact scores are incorporated to assess the predicted biological consequences of genomic variants. These include established deleteriousness predictors such as CADD,

SIFT, and PolyPhen.

Rather than treating individual scores in isolation, composite indicators are constructed to classify variants into interpretable impact categories such as likely benign, uncertain, or potentially deleterious. These indicators are explicitly designed to remain transparent, auditable, and suitable for downstream feature inspection.

10.3 Population Frequency Integration and Leakage Prevention

Population allele frequencies are integrated using aggregated population datasets to distinguish rare variants from common polymorphisms. Variants are classified into rare, low-frequency, and common categories based on population prevalence thresholds.

These classifications serve both analytical and methodological purposes. From a modeling perspective, they reduce information leakage by preventing downstream models from implicitly learning population prevalence as a proxy for pathogenicity without explicit control.

10.4 Feature Grouping and Confidence Scoring

Variant-level confidence metrics are computed by combining conservation strength, functional impact indicators, population rarity, and clinical review quality signals. These grouped features provide monotonic, interpretable confidence scores that reflect cumulative biological evidence.

The resulting confidence metrics are designed to be stable across reprocessing cycles and suitable for use in downstream analytics and machine learning workflows without introducing opaque heuristics or undocumented transformations.

Formal feature lineage principles and traceability rules are documented in Appendix D.

Chapter 11

Classical Machine Learning

This chapter corresponds to Phase 7 of the DNA Gene Mapping Project and focuses on transforming enriched biological data into structured, machine-learning-ready feature sets aligned with real clinical research and pharmaceutical analytics use cases. The emphasis of this phase is on feature construction, scoring logic, and interpretability rather than large-scale model deployment.

11.1 Disease-Level Feature Construction

Classical machine learning readiness in this project is achieved by aggregating curated Gold-layer features across variants, genes, and diseases into structured, interpretable feature sets.

Disease-centric features are constructed by aggregating variant-level and gene-level relationships for each normalized disease entity. Derived metrics include disease variant counts, pathogenic variant ratios, gene diversity measures, and disease complexity indicators.

These features support downstream polygenic risk analysis, disease stratification, and prioritization workflows without embedding model-specific assumptions.

Figure 11.1 illustrates how Gold-level biological features are combined, scored, and transformed into machine-learning-ready inputs suitable for downstream analytical workflows.

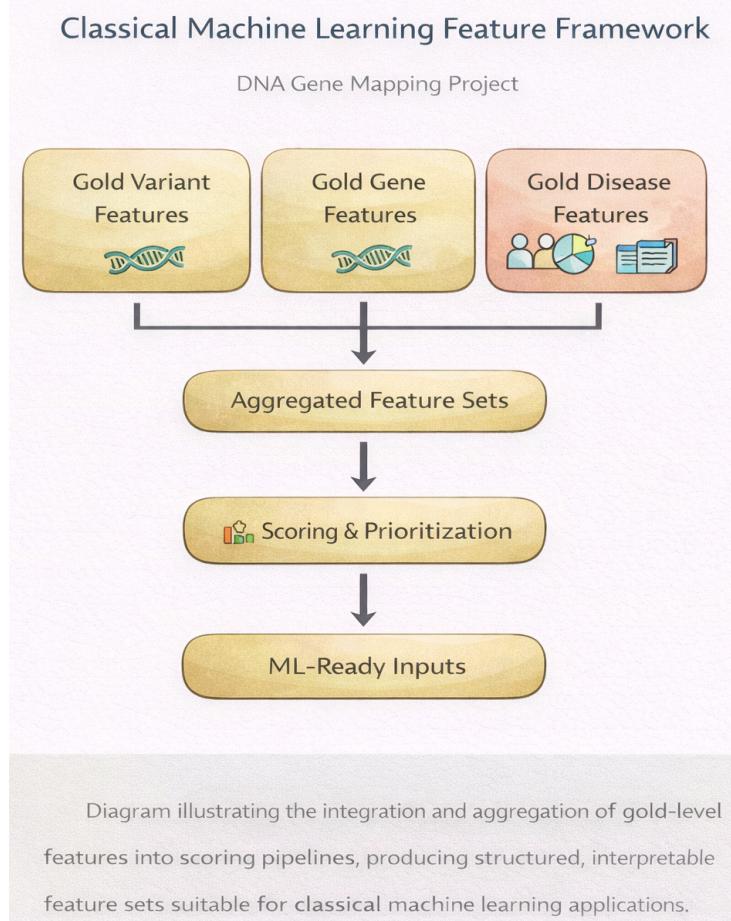


Figure 11.1: Classical machine learning feature framework showing the aggregation of Gold-layer variant, gene, and disease features into structured feature sets, followed by scoring and prioritization to produce interpretable, machine-learning-ready inputs.

11.2 Gene Prioritization Feature Framework

Gene prioritization features are derived by quantifying how many genes and variants contribute to a given disease. Diseases are categorized into monogenic, oligogenic, and polygenic classes based on gene participation counts.

At the variant level, polygenic contribution indicators are computed to represent relative influence within genetically complex diseases. These features form the foundation for gene ranking and prioritization in clinical research and discovery pipelines.

11.3 Gene–Disease Association Scoring

Genes are evaluated for clinical and analytical relevance using aggregated signals including variant burden, pathogenic variant counts, disease diversity, and annotation depth.

Composite gene utility scores are computed to support prioritization workflows com-

monly used in diagnostic research, gene panel design, and translational analytics. These scores are designed to be interpretable and auditable rather than optimized for opaque predictive performance.

11.4 Evaluation and Validation Signals

Evaluation signals are introduced to assess the stability and biological plausibility of constructed features. Disease burden indicators are used to identify diseases with disproportionately high pathogenic variant loads or unusually broad gene involvement.

These signals enable identification of genetically complex, clinically severe, or under-characterized diseases, supporting both exploratory research and downstream modeling readiness.

11.5 Explainability and Transparency

Explainability is treated as a first-class design requirement throughout the classical machine learning framework. All derived features and scores are constructed using explicit, documented logic and traceable source inputs.

Rather than relying on black-box predictions, the system emphasizes transparent feature definitions, monotonic scoring behavior, and clear biological interpretation, ensuring that downstream users can understand, audit, and trust analytical outputs.

Model feature traceability and explainability constraints align with the lineage framework defined in Appendix D.

Chapter 12

Deep Learning and Graph Modeling

This chapter corresponds to Phase 8 of the DNA Gene Mapping Project and focuses on preparing the dataset, feature representations, and architectural foundations for advanced machine learning techniques. The intent of this phase is to enable future deep learning and graph-based modeling without prematurely committing to resource-intensive implementations or exceeding current computational constraints.

12.1 Why Deep Learning Is Optional but Valuable

Not all analytical use cases in genomics require deep learning. Many disease prioritization, gene ranking, and variant triage tasks can be addressed effectively using classical machine learning and interpretable scoring frameworks.

However, deep learning offers additional value in capturing non-linear interactions, higher-order dependencies, and latent biological patterns that are difficult to encode explicitly. The system is therefore designed to support deep learning as an optional analytical extension rather than a mandatory dependency.

The conceptual role of deep learning within the system is defined in Appendix J.

12.2 Model-Agnostic Feature Preparation

All engineered features are designed to be model-agnostic. Clear separation is maintained between identifiers, clinical labels, explanatory features, and derived scores to prevent data leakage and unintended label encoding.

This design enables the same feature sets to be reused across classical machine learning, deep learning, and graph-based models without restructuring the underlying data pipeline.

12.3 Convolutional Neural Networks for Sequence-Based Modeling

Convolutional Neural Networks (CNNs) are identified as a suitable future modeling approach for learning localized patterns in biological sequences and structured feature matrices. Potential applications include motif detection in genomic regions, localized conservation pattern analysis, and protein functional region characterization.

CNN-based models are treated as optional, experimental extensions. Feature preparation is designed to allow CNN inputs to be derived from curated tables without embedding sequence processing logic directly into the core pipeline.

12.4 Graph Neural Networks for Biological Networks

Graph Neural Networks (GNNs), including Graph Convolutional Networks, are identified as a high-value future direction for modeling biological networks involving genes, proteins, variants, and diseases.

In this context, nodes represent biological entities while edges capture curated relationships such as gene–disease associations, protein–protein interactions, and variant–gene links. The current system explicitly prepares graph-compatible representations while deferring full graph-based model training to later phases.

12.5 Transformer-Based Contextual Modeling

Transformer architectures are identified as a potential future approach for modeling long-range dependencies and contextual relationships within biological data. Possible applications include protein sequence dependency modeling, ordered variant context analysis, and disease progression signal extraction.

Given current execution constraints, transformer-based models are not implemented directly. However, feature representations and data organization are designed to remain compatible with future transformer-based experimentation.

12.6 Bias Awareness and Mitigation Strategies

Potential sources of bias are explicitly identified, including population representation imbalance, reporting bias in clinical databases, and annotation incompleteness across genes and diseases.

Mitigation strategies are documented and applied where feasible, including feature normalization, confidence weighting based on evidence strength, and exclusion of unstable or sparsely supported signals from training targets.

12.7 Execution Constraints and Architectural Planning

Given current hardware and execution constraints, deep learning workflows are planned as controlled, experimental extensions rather than fully deployed systems. Execution is designed to leverage Databricks Community Edition for data preparation and limited model experimentation, without reliance on GPUs or managed cloud services.

This approach allows architectural validation and feature readiness testing while remaining compatible with local-first execution principles.

12.8 What Is Safe to Implement Now vs Later

Certain deep learning components are considered safe to implement incrementally, including feature extraction pipelines, embedding generation from intermediate model layers, and biological graph construction logic.

More resource-intensive activities—such as large-scale neural network training, extensive hyperparameter optimization, or end-to-end deep learning pipelines—are intentionally deferred. This phased approach ensures that scientific validity and system integrity are preserved without overextending computational resources.

Execution-level deep learning constraints and pipelines are described in Appendix K.

Chapter 13

Power BI Analytics

This chapter corresponds to Phase 9 of the DNA Gene Mapping Project and focuses on the analytical and visualization layer built on top of the Gold data model. The objective of this phase is to enable exploratory analysis, domain-specific reporting, and interpretability of genomic insights using local, desktop-based business intelligence tooling.

13.1 Dashboard Architecture

The analytics layer is implemented using Power BI Desktop in local execution mode. Gold-layer feature tables serve as the primary data source for all dashboards, ensuring that visualizations are built on curated, stable, and analytically meaningful datasets.

Dashboards are designed using a modular architecture, with each analytical use case implemented as an independent report page. This approach allows individual dashboards to evolve without impacting other analytical views while preserving consistent data definitions and metrics.

13.2 Gene Analytics Dashboards

Gene-focused dashboards provide interactive views into gene-level annotations, disease associations, and prioritization scores. Visual elements include gene-level feature summaries, disease diversity indicators, and prioritization rankings derived from Gold-layer features.

These dashboards support exploratory research workflows, allowing users to identify genes of interest based on functional relevance, disease burden, and analytical confidence signals.

13.3 Variant Analytics Dashboards

Variant dashboards focus on the exploration of clinically and biologically relevant variants. Visualizations include variant classification distributions, pathogenicity confidence indicators, population frequency stratification, and gene-level aggregation views.

These dashboards are designed to support variant triage, exploratory analysis, and hypothesis generation rather than automated clinical decision-making.

13.4 Disease Risk and Complexity Dashboards

Disease-centric dashboards present aggregated disease-level features, including variant burden metrics, polygenic complexity indicators, and gene participation diversity.

These views enable comparative analysis across diseases, identification of genetically complex conditions, and prioritization of diseases for further research or modeling.

13.5 Local Execution and Usage Constraints

All dashboards are executed locally using Power BI Desktop. No Power BI Service, cloud gateways, or online publishing features are used.

Data ingestion into Power BI is performed via flat-file exports or direct connections to local analytical stores where applicable. This design ensures offline operability, cost-free execution, and full compliance with the project's local-first architectural constraints.

Chapter 14

Advanced Analytics and Visualization Layer

This chapter focuses on transforming curated analytical data into interpretable, decision-ready insights through structured analytics and visualization. The goal of this layer is to bridge complex genomic feature engineering with stakeholder-facing outputs suitable for researchers, clinical analysts, and domain experts.

Unlike upstream engineering and modeling components, this layer emphasizes interpretability, transparency, and exploratory analysis rather than algorithmic complexity.

14.1 Analytical Objectives

The primary objective of the analytics layer is to expose gene, variant, and disease intelligence through stable, explainable metrics derived from Gold-layer tables.

Key analytical questions addressed include:

- Which genes contribute most significantly to a disease
- Which variants drive pathogenic burden
- How diseases differ in genetic complexity
- Which genes should be prioritized for further investigation

All analytical outputs are derived exclusively from curated Gold-layer features to ensure consistency, traceability, and reproducibility.

14.2 Visualization Consumption Strategy

Analytical outputs are designed to be consumed by visualization and reporting tools without requiring direct access to raw or intermediate data. This separation ensures that

analytical logic remains centralized and auditable.

Visualization tools operate strictly in read-only mode and do not perform transformations that would alter analytical semantics or feature definitions.

14.3 Dashboard Design Themes

Dashboards are organized around thematic analytical views rather than raw table exposure. Core dashboard themes include:

Disease Overview Dashboards Present disease complexity, pathogenic burden, gene diversity, and polygenic indicators.

Gene Prioritization Dashboards Rank genes based on clinical utility, pathogenic variant burden, disease associations, and annotation depth.

Variant Landscape Dashboards Visualize variant distributions by clinical significance, population frequency, conservation strength, and disease involvement.

Each dashboard supports filtering by disease, gene, chromosome, and clinical relevance attributes.

14.4 Data Refresh and Governance

Data refresh cycles are explicit, controlled, and versioned. No live mutation of analytical data is permitted within the visualization layer.

All metrics and visual outputs are traceable back to Gold-layer tables, ensuring auditability, reproducibility, and regulatory defensibility.

Chapter 15

Machine Learning Model Implementation and Evaluation

This chapter transitions from feature readiness into the implementation and evaluation of machine learning models. The focus is on validating the analytical usefulness of engineered features rather than building production-scale predictive systems.

Modeling activities are executed under controlled conditions to preserve interpretability, reproducibility, and biological plausibility.

15.1 Execution Environment

All machine learning workloads are executed within Databricks Community Edition to avoid reliance on local hardware for computationally intensive operations.

The environment supports:

- Scalable feature loading from Gold-layer tables
- Reproducible model training workflows
- Deterministic experiment execution

No managed cloud services, GPU acceleration, or proprietary ML platforms are used.

15.2 Initial Modeling Use Cases

Initial modeling efforts focus on well-defined, high-confidence objectives:

Variant Pathogenicity Classification Classification models are trained using curated pathogenic labels combined with conservation, functional impact, and review-quality features.

Gene Prioritization Modeling Regression and ranking models score genes based on clinical utility, disease diversity, and variant burden.

Disease Complexity Prediction Classification models distinguish monogenic, oligogenic, and polygenic diseases using aggregated disease-level features.

15.3 Model Selection Philosophy

Early models prioritize interpretability and stability over predictive complexity. Preferred model families include:

- Regularized linear models
- Gradient-boosted decision trees

Deep learning models are intentionally deferred until feature distributions, biases, and stability characteristics are well understood.

15.4 Evaluation Strategy

Evaluation focuses on biologically and clinically meaningful metrics rather than purely statistical performance.

Metrics include:

- Precision and recall for pathogenic variant classification
- Ranking consistency for gene prioritization
- Stability across disease subgroups

Cross-validation strategies are designed to prevent gene-level and disease-level leakage across folds.

Chapter 16

System Architecture and Application Layer Design

This chapter defines the system architecture responsible for delivering analytical and modeling outputs through programmatic and interactive interfaces. The application layer is designed to expose curated intelligence while preserving data integrity and access control.

The architecture emphasizes modularity, separation of concerns, and controlled access to analytical assets.

The logical architecture layers and technology boundaries referenced in this chapter are formally defined in Appendix E.

16.1 Backend Architecture

The backend is designed as a modular, service-oriented system responsible for exposing read-only analytical endpoints.

Core backend components include:

- FastAPI for API services
- Databricks for analytical computation
- PostgreSQL for structured storage and integration (future state)

APIs are intentionally restricted to read-only operations to prevent unauthorized mutation of curated analytical data.

16.2 Frontend Architecture

The frontend consumes backend APIs and presents analytical outputs through intuitive, workflow-oriented interfaces.

The planned frontend stack includes:

- React for component-based UI development
- TypeScript for type safety and maintainability
- REST-based data access

Frontend components are designed around user workflows such as gene exploration, disease analysis, and variant triage rather than raw data exposure.

16.3 Project Structure and Domain Separation

The overall project structure is organized into clearly separated domains, including:

- Data ingestion and processing
- Analytics and feature engineering
- Machine learning and modeling
- API services
- Frontend applications
- Documentation

This structure enforces clear ownership boundaries and supports long-term maintainability.

16.4 Security and Access Control

Security is enforced through role-based access control at both the API and data layers.

No direct database access is exposed to frontend clients. All interactions are mediated through validated service endpoints, ensuring controlled access, traceability, and auditability.

Chapter 17

Backend API Layer

This chapter defines the backend application layer responsible for exposing curated analytical data and machine learning outputs through programmatic interfaces. The backend serves as a controlled access boundary between internal analytical assets and downstream consumers such as web applications, dashboards, and external analytical tools.

The backend is explicitly designed for read-only access to prevent unauthorized mutation of curated biological data, analytical features, or model-derived outputs.

17.1 FastAPI Architecture

The backend is implemented using FastAPI due to its lightweight design, strong typing support, and suitability for analytical API services.

The FastAPI application is organized into modular components that separate routing, data access, schema validation, and business logic. This structure ensures that analytical logic remains decoupled from API presentation and that changes in upstream data models do not directly impact API consumers.

FastAPI's native support for OpenAPI specifications enables automatic API documentation, improving transparency and ease of integration for downstream users.

17.2 API Endpoint Design

API endpoints are designed around analytical use cases rather than raw database tables. Each endpoint exposes a stable, well-defined contract that maps directly to curated Gold-layer tables or derived analytical views.

Primary endpoint categories include:

- Gene-centric analytical endpoints
- Variant-centric analytical endpoints

- Disease-centric analytical endpoints
- Aggregated summary and ranking endpoints

Endpoints support controlled filtering, pagination, and parameterized queries while explicitly avoiding free-form query execution. This approach preserves performance predictability and prevents misuse of analytical resources.

17.3 Security Considerations

Security is enforced through a combination of architectural constraints and explicit access controls. The backend exposes read-only endpoints only and does not provide any mutation, ingestion, or administrative operations.

Access control is designed to support role-based authorization in future extensions. Sensitive internal identifiers and raw source metadata are not exposed directly through API responses unless explicitly required for analytical transparency.

All external access to analytical data is mediated through validated API schemas, ensuring consistent data representation and reducing the risk of schema drift or accidental data leakage.

17.4 Model Serving Strategy

Machine learning models are not exposed as black-box prediction services. Instead, model outputs are materialized as features within Gold-layer tables and served through the same analytical APIs as other curated metrics.

This strategy ensures:

- Interpretability of model-derived outputs
- Consistent versioning and reproducibility
- Unified access patterns for analytical and ML features

Where necessary, lightweight inference endpoints may expose precomputed scores or rankings, but real-time model execution is intentionally avoided at this stage to preserve system stability and execution feasibility.

The backend therefore functions as a delivery layer for analytical intelligence rather than an online prediction engine.

Constraints governing local execution and model serving are defined in Appendix N.

Chapter 18

Web Application Frontend

This chapter defines the web application frontend responsible for presenting analytical insights to end users through an interactive, intuitive interface. The frontend serves as the primary human-facing layer of the system, enabling exploration of genes, variants, and diseases without exposing raw data structures or backend complexity.

The frontend is designed to consume backend APIs exclusively and does not perform independent analytical transformations.

18.1 React Architecture

The frontend is implemented using a React-based architecture to support modular, component-driven development and responsive user interfaces.

React is selected for its ecosystem maturity, strong community support, and suitability for building analytical dashboards that require dynamic data loading, filtering, and state management. TypeScript is used throughout the frontend to enforce type safety and reduce integration errors between frontend components and backend APIs.

The application follows a unidirectional data flow, ensuring predictable state transitions and simplified debugging.

18.2 Component Structure

Frontend components are organized by analytical domain rather than technical function. This structure aligns UI components directly with user workflows and analytical use cases.

Core component categories include:

- Gene exploration components
- Variant analysis components
- Disease overview and comparison components

- Shared visualization and filtering components

Reusable components such as tables, charts, filters, and pagination controls are centralized to avoid duplication and ensure consistent interaction patterns across the application.

18.3 API Integration

All data displayed in the frontend is retrieved via validated backend API endpoints. The frontend does not directly access databases or analytical storage layers.

API integration is designed to support:

- Parameterized data queries
- Controlled pagination for large result sets
- Explicit error handling and fallback states

Data fetching logic is isolated from presentation components to preserve separation of concerns and simplify future refactoring or backend evolution.

The local analytics and visualization consumption model is described in Appendix I.

18.4 User Experience Principles for Genomics

The frontend user experience is designed with domain-specific considerations for genomics and biomedical analytics.

Key UX principles include:

- Progressive disclosure of complex information to avoid cognitive overload
- Clear labeling of analytical confidence and evidence strength
- Consistent use of biological terminology aligned with curated schemas
- Avoidance of implicit clinical recommendations or decision cues

Visual elements emphasize interpretability and context rather than aesthetic complexity. The interface is intended to support exploratory analysis, hypothesis generation, and educational use rather than clinical decision-making.

The frontend therefore functions as a transparent analytical lens into the system rather than an automated recommendation engine.

Chapter 19

System Integration and Deployment

This chapter describes how the analytical, machine learning, backend, and frontend components of the DNA Gene Mapping Project are integrated into a cohesive system. The focus is on end-to-end data flow, controlled execution, and future extensibility rather than production-scale deployment.

System integration is designed to preserve data integrity, analytical traceability, and execution feasibility under local-first constraints.

19.1 End-to-End System Flow

The integrated system follows a deterministic, layered execution flow beginning with curated analytical data and ending with user-facing insights.

At a high level:

- Curated Silver- and Gold-layer tables serve as the analytical foundation
- Machine learning outputs are materialized as Gold-layer features
- Backend APIs expose read-only analytical endpoints
- Frontend applications consume APIs for interactive exploration
- Visualization tools access the same curated outputs for reporting

This flow ensures that all downstream consumption is derived from validated, versioned analytical assets rather than raw or transient data.

The canonical project folder structure and codebase organization are documented in Appendix O.

19.2 Execution and Orchestration Strategy

System execution is orchestrated through controlled, repeatable workflows rather than automated production pipelines. Data processing, feature generation, and model execution are triggered explicitly to maintain transparency and reproducibility.

Long-running analytical tasks are executed within Databricks Community Edition, while application-layer components operate independently and consume only finalized outputs.

This strategy avoids tight coupling between computation-heavy processes and interactive services, improving system stability.

19.3 Monitoring and Observability

Monitoring focuses on correctness, data completeness, and execution integrity rather than real-time performance metrics.

Key observability signals include:

- Successful completion of analytical and ML workflows
- Schema validation and data consistency checks
- API availability and contract compliance

Failures or inconsistencies are surfaced through explicit logging and validation reports rather than automated remediation, supporting manual inspection and auditability.

19.4 Scalability Considerations

Scalability is approached as a controlled, incremental concern rather than a primary design driver. The system is structured to scale analytically through distributed computation while keeping application components lightweight.

Design choices that support scalability include:

- Schema-driven data processing
- Layered separation of analytical and application concerns
- Deferred reliance on persistent relational storage

These decisions ensure that scaling efforts can be introduced gradually without requiring architectural redesign.

19.5 Future Portability

Although currently executed in a local-first environment, the system is architected to remain portable across execution contexts.

Portability considerations include:

- Clear separation between computation and presentation layers
- Explicit data contracts between system components
- Avoidance of vendor-specific service dependencies

This design allows future migration to alternative execution environments or infrastructure models while preserving analytical logic and system semantics.

Chapter 20

ML Feature Inventory (High-Level)

This chapter provides a consolidated, high-level inventory of machine learning features generated by the DNA Gene Mapping Project. The objective of this chapter is to document the scope and intent of features made available for analytical modeling without enumerating low-level schema details.

All features described in this chapter are derived from curated Silver- and Gold-layer tables and are designed to be interpretable, reproducible, and suitable for downstream machine learning and analytical workflows.

Core biological entity features are defined in Appendix B, while this chapter focuses on derived and analytical feature groupings.

Derived analytical and machine learning features, including disease association, polygenic risk, and gene prioritization metrics, are documented in Appendix C.

The lineage, traceability, and enforcement rules governing all analytical and machine learning features are documented in Appendix D.

20.1 Gene-Level Features

Gene-level features capture functional relevance, disease involvement, and analytical utility at the gene entity level. These features support gene prioritization, disease association modeling, and target discovery workflows.

Key gene feature categories include:

- Gene identity and normalization indicators
- Functional annotation summaries
- Disease association counts and diversity metrics
- Pathogenic variant burden indicators
- Clinical relevance and prioritization scores

Gene features are designed to remain stable across modeling iterations and to support both ranking-based and classification-based analytical tasks.

20.2 Variant-Level Features

Variant-level features represent the core analytical signals used for pathogenicity assessment, prioritization, and disease impact analysis.

Key variant feature categories include:

- Clinical significance classifications
- Review status and evidence strength indicators
- Population frequency and rarity flags
- Conservation and functional impact scores
- Aggregated confidence and quality metrics

Variant features are explicitly designed to avoid label leakage by separating explanatory features from clinical outcome labels used in supervised modeling.

20.3 Disease-Level Features

Disease-level features aggregate gene and variant information into disease-centric analytical representations. These features enable modeling of disease complexity, polygenicity, and genetic burden.

Key disease feature categories include:

- Gene participation counts per disease
- Pathogenic variant burden metrics
- Polygenic complexity indicators
- Disease-level confidence and stability scores

Disease features are intended to support disease stratification, comparative analysis, and risk modeling rather than direct clinical prediction.

20.4 Cross-Entity and Derived Features

In addition to entity-specific features, the system generates derived features that capture relationships across genes, variants, and diseases.

Examples include:

- Gene–disease association strength indicators
- Variant contribution scores within polygenic diseases
- Network-based aggregation metrics

These features enable more expressive modeling while preserving interpretability and traceability to underlying biological entities.

20.5 Feature Governance Principles

All machine learning features adhere to strict governance principles:

- Every feature is traceable to a source table and transformation logic
- Feature definitions are versioned and documented
- No feature is introduced without a defined analytical purpose
- Feature semantics remain consistent across modeling iterations

This governance framework ensures that machine learning experimentation remains scientifically grounded, auditable, and reproducible.

Chapter 21

Feature Lineage Documentation

This chapter documents how analytical and machine learning features in the DNA Gene Mapping Project are derived, traced, and governed across the data lifecycle. Feature lineage is treated as a first-class design requirement to ensure transparency, reproducibility, and long-term maintainability.

Rather than relying on external lineage tooling, lineage is enforced through explicit schema design, deterministic transformation logic, and structured documentation.

21.1 Source-to-Feature Traceability

Every feature generated in the system can be traced back to its authoritative source datasets. Lineage is preserved across all transformation stages, including ingestion, normalization, enrichment, and feature engineering.

For each feature, lineage documentation includes:

- Upstream source tables and identifiers
- Transformation logic applied
- Intermediate representations (if any)
- Final analytical or modeling usage

This traceability ensures that analytical results can be audited and biologically interpreted rather than treated as opaque outputs.

21.2 Lineage Across Data Layers

Lineage is maintained consistently across all logical data layers.

Silver Layer Features in the Silver layer represent normalized and biologically enriched entities. Lineage at this stage focuses on correctness, identifier consistency, and preservation of source semantics.

Gold Layer Gold-layer features are explicitly engineered for analytics and machine learning. Lineage documentation at this stage emphasizes feature intent, aggregation logic, and analytical meaning rather than raw biological annotation.

This layered lineage approach prevents downstream consumers from relying on unstable or ambiguously derived features.

21.3 Auditing and Reproducibility

Feature lineage enables deterministic re-creation of analytical outputs. Given a specific dataset version and transformation sequence, features can be regenerated without ambiguity.

Auditing is supported through:

- Versioned schemas
- Explicit transformation notebooks or scripts
- Controlled execution order

This design ensures that results remain reproducible across time, environments, and future extensions of the system.

21.4 Explainability by Design

Explainability is embedded at feature design time rather than retrofitted after modeling. Each feature is introduced with a clear biological or analytical rationale and an intended downstream use.

Model-derived outputs are materialized as interpretable features rather than exposed as black-box predictions. This preserves trust, supports validation, and aligns with the project's research-oriented and educational objectives.

Feature lineage therefore functions not only as a technical mechanism, but as a conceptual bridge between biological knowledge, analytical reasoning, and machine learning experimentation.

Chapter 22

Conclusion

The DNA Gene Mapping Project demonstrates how complex genomic and biomedical data can be transformed into a structured, analytically robust, and machine-learning-ready system under realistic constraints. The project moves beyond isolated bioinformatics scripts to present an end-to-end platform that integrates data engineering, analytics, modeling readiness, and system design.

By emphasizing schema discipline, feature lineage, and layered architecture, the system ensures correctness, traceability, and long-term extensibility. Authoritative public datasets are unified into coherent gene, variant, disease, and protein representations that support advanced analytical use cases without sacrificing interpretability.

The project intentionally adopts a local-first execution strategy, proving that serious genomics analytics and modeling readiness can be achieved without reliance on enterprise cloud infrastructure. Constraints related to hardware, storage, and compute are treated as architectural drivers rather than limitations, resulting in disciplined and efficient system design.

Equally important, the platform establishes clear boundaries around intended use, ethical considerations, and analytical responsibility. Outputs are designed for research, education, and exploratory analysis, not for clinical diagnosis or decision-making.

Taken together, this work provides a realistic blueprint for building genomics analytics platforms that balance scientific rigor, engineering discipline, and practical feasibility. It serves as a strong foundation for future research, system evolution, and professional development in biomedical data science and genomics engineering.

Acknowledgement of Data Sources

This project uses publicly available biomedical and genomic data resources, including ClinVar, MedGen, OMIM, RefSeq, UniProt, gnomAD, dbVar, and the Genetic Testing Registry (GTR). These resources are maintained by international research institutions, public databases, and scientific consortia and are widely used within the genomics and biomedical research communities.

All data sources are accessed according to their respective licenses, terms of use, and citation requirements. No proprietary, restricted-access, or patient-identifiable data is included at any stage of this project. Datasets are used strictly for research, analytical, and educational purposes.

The author acknowledges the substantial scientific effort involved in the curation, validation, and long-term maintenance of these resources. The availability of such high-quality public biomedical datasets enables reproducible research, transparent analysis, and independent innovation in genomics, disease modeling, and variant interpretation.

References

Appendices

Appendix A

Data Layer Schemas

This appendix documents the logical structure and responsibilities of the data layers used throughout the DNA Gene Mapping Project. The system follows a layered schema model that separates authoritative reference data, normalized biological entities, and derived analytical features to ensure correctness, traceability, and analytical stability across the entire pipeline.

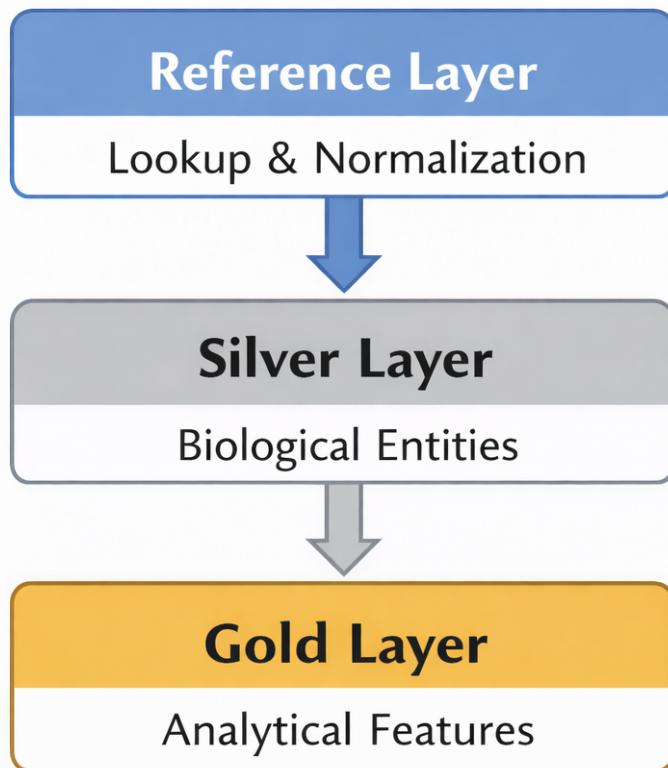


Figure A.1: Logical relationship between reference, silver, and gold data layers. The reference layer provides authoritative normalization and controlled vocabularies, the silver layer models normalized biological entities, and the gold layer contains analytical and machine-learning-ready feature tables derived from silver-layer data.

A.1 Reference Layer Table Schemas

The reference layer contains authoritative lookup datasets used across the entire pipeline. These tables are treated as immutable sources of truth and are never modified by downstream analytical processes.

The primary purposes of the reference layer include identifier normalization, cross-database alignment, controlled vocabulary resolution, and disease and ontology mapping. Typical entities stored in this layer include MedGen concepts, OMIM identifiers, MONDO disease mappings, and gene identifier crosswalks.

Each reference table uses stable identifiers and avoids derived or heuristic logic. Table A.1 summarizes the core reference-layer tables and their primary identifiers.

Table A.1: Reference layer core tables and primary identifiers

Table Name	Primary Key	Description
reference.medgen_concepts	medgen_id	Canonical disease concepts and preferred terms
reference.omim_mappings	omim_id	OMIM identifiers aligned to MedGen concepts
reference.mondo_mappings	mondo_id	MONDO ontology cross-references
reference.gene_id_crosswalk	entrez_gene_id	Gene identifier normalization across databases

A.2 Silver Layer Table Schemas

The silver layer represents cleaned, normalized, and biologically enriched entities. This layer forms the analytical backbone of the project and serves as the authoritative biological representation used by all downstream analytics.

Key characteristics of the silver layer include one biological entity per table, strong primary keys, fully normalized identifiers, and deterministic, reproducible transformations. No machine-learning features or predictive scores are stored in this layer.

Table A.2 lists the core silver-layer tables and the biological entities they model.

Table A.2: Silver layer core biological entity tables

Table Name	Primary Key	Entity Modeled
silver.genes	gene_id	Gene entity
silver.variants	variant_id	Genetic variant
silver.diseases	medgen_id	Disease entity
silver.gene_disease_links	link_id	Gene–disease association
silver.proteins	protein_id	Protein entity
silver.transcripts	transcript_id	Transcript entity

Silver-layer data is treated as biologically authoritative within the system and is designed to remain stable across analytical and modeling iterations.

A.3 Gold Layer Table Schemas

The gold layer contains analytical and machine-learning-ready feature tables derived from silver-layer entities. These tables are denormalized for analytical efficiency and optimized for consumption by statistical models, machine learning pipelines, and visualization tools.

Gold-layer schemas are feature-driven and explicitly versioned to preserve feature lineage and reproducibility. Gold tables are not treated as authoritative biological records; instead, they represent analytical constructs derived from curated biological data.

Table A.3 summarizes the primary gold-layer feature tables and their analytical focus.

Table A.3: Gold layer analytical and machine learning feature tables

Table Name	Primary Key	Feature Focus
gold.variant_features	variant_id	Variant-level analytical and ML features
gold.gene_features	gene_id	Gene prioritization and functional scores
gold.disease_features	disease_id	Disease complexity and risk features
gold.gene_prioritization	gene_id	Composite gene ranking metrics

The layered schema design ensures that raw biological truth remains isolated from analytical abstraction, enabling stable analytics, transparent feature lineage, and controlled system evolution.

Appendix B

Feature Dictionary (Core Biological Entities)

This appendix defines the core biological features modeled within the DNA Gene Mapping Project. These features represent normalized, biologically meaningful attributes associated with genes, variants, diseases, and their relationships. The feature definitions provided here apply to the curated silver layer and selected gold-layer entities and are intended to support interpretability, auditability, and downstream analytical usage.

B.1 Gene Features (Silver Layer)

This section documents the primary features associated with gene entities. Each gene is represented as a fully enriched biological object serving as a central integration point for variant, disease, protein, and functional annotations.

The primary identifier for this entity is `gene_id`. Table B.1 summarizes the key gene-level features captured in the silver layer.

Gene features act as the primary biological integration hub and are referenced extensively by downstream analytical and machine learning workflows.

B.2 Variant Features (Silver Layer)

This section defines features associated with normalized genetic variants. Variant entities represent individual genomic alterations with curated biological and clinical context.

The primary identifier for this entity is `variant_id`. Table B.2 summarizes the core variant-level features.

Variant features support both clinical analytics and downstream machine learning feature generation.

Table B.1: Core gene features (silver layer)

Feature Category	Description
Gene symbols and aliases	Official gene symbol and curated alternative names resolved across source databases
Functional classifications	High-level functional roles inferred from curated annotations
Druggability indicators	Flags indicating known or potential drug target relevance
Protein function categories	Functional groupings derived from linked protein annotations
Genomic location	Chromosome, genomic coordinates, and cytogenetic context
Pharmacogenomic annotations	Indicators of relevance to drug metabolism and response

Table B.2: Core variant features (silver layer)

Feature Category	Description
Allele mappings	Canonical allele identifiers and relationships across variant sources
Transcript-level annotations	Variant impact relative to annotated transcripts
Protein-level changes	Amino-acid level consequences derived from curated mappings
Clinical significance	Standardized clinical interpretation categories
Disease associations	Linked disease concepts derived from curated sources
Variant classification flags	Boolean indicators supporting efficient filtering and prioritization

B.3 Disease Features (Silver Layer)

This section documents features associated with normalized disease entities. Diseases are treated as first-class analytical objects rather than free-text labels.

The primary identifier for this entity is `medgen_id`. Table B.3 summarizes the disease-level features.

Table B.3: Core disease features (silver layer)

Feature Category	Description
Disease names	Preferred disease terms enforced across the system
Ontology alignments	Cross-references to external disease ontologies
Source provenance	Originating databases contributing disease definitions
Disease category indicators	High-level disease classification flags

Disease features enable consistent analytical grouping, aggregation, and hierarchical analysis.

B.4 Gene–Disease Association Features (Silver Layer)

This section defines explicit gene–disease relationship features. These associations capture curated biological evidence linking genes to diseases.

Gene–disease relationships are modeled as separate analytical entities to support many-to-many mappings. Table B.4 summarizes the association features.

Table B.4: Gene–disease association features (silver layer)

Feature Category	Description
Gene identifier	Reference to the associated gene entity
Disease identifier	Reference to the associated disease entity
Association source	Database or resource supporting the association
Evidence type	Qualitative classification of supporting evidence

These features enable downstream network-based analysis and polygenic disease modeling.

B.5 Disease Machine Learning Features (Gold Layer)

This section documents disease-level features engineered specifically for analytical modeling and machine learning. These features are derived from curated silver-layer entities

and stored in gold-layer tables.

Table B.5 summarizes the primary disease-level machine learning feature groups.

Table B.5: Disease-level machine learning features (gold layer)

Feature Group	Description
Disease complexity metrics	Indicators of genetic heterogeneity and polygenicity
Polygenic risk indicators	Measures of multi-gene contribution to disease risk
Pathogenic burden scores	Aggregated pathogenic variant impact metrics
Gene prioritization metrics	Scores supporting disease-centric gene ranking

These features are designed to be directly consumable by machine learning pipelines without additional joins or transformations.

Appendix C

Feature Dictionary (Derived & Analytical Features)

This appendix documents derived and analytical features engineered from curated biological entities. Unlike core entity features, which describe biological attributes directly, the features defined here represent analytical abstractions designed to support disease modeling, gene prioritization, and machine learning workflows.

All features described in this appendix are derived deterministically from silver-layer entities and persisted in gold-layer analytical tables. These features are explicitly versioned and governed to preserve interpretability and reproducibility.

C.1 Disease Association Features

Disease association features quantify the breadth and strength of relationships between genes, variants, and diseases. These features enable analytical characterization of disease relevance and genetic burden without embedding causal assumptions.

Table C.1 summarizes the primary disease association features.

Table C.1: Derived disease association features

Feature Name	Description
disease_count	Number of distinct diseases associated with a gene
has_disease_association	Boolean indicator of any disease linkage
disease_association_strength	Categorical measure of association strength based on disease count

These features support downstream disease prioritization and complexity analysis.

C.2 Polygenic Risk Features

Polygenic risk features characterize diseases based on the number and diversity of contributing genes and variants. These features enable differentiation between monogenic and complex diseases and support polygenic risk modeling.

Table C.2 summarizes the core polygenic risk features.

Table C.2: Derived polygenic risk features

Feature Name	Description
disease_gene_count	Number of genes contributing to a disease
is_polygenic_disease	Boolean indicator for multi-gene diseases
disease_pathogenic_ratio	Fraction of pathogenic variants associated with a disease
polygenic_risk_contribution	Per-variant contribution score within a polygenic context

These features are designed to support disease stratification, genetic complexity analysis, and downstream machine learning models.

C.3 Gene Prioritization Features

Gene prioritization features aggregate biological, clinical, and analytical signals to support systematic ranking of genes for research and translational investigation.

Table C.3 summarizes the primary gene prioritization features.

Table C.3: Derived gene prioritization features

Feature Name	Description
gene_clinical_utility_score	Composite score reflecting pathogenicity, evidence strength, and disease diversity
gene_priority_tier	Tiered classification of gene relevance for clinical and research use
is_clinically_actionable	Boolean flag indicating high-confidence actionable genes

Gene prioritization features are intentionally transparent and interpretable to support trust in analytical outputs and downstream decision-making.

Appendix D

Feature Lineage Documentation

This appendix documents the feature lineage framework used throughout the DNA Gene Mapping Project. Feature lineage is treated as a first-class design requirement to ensure transparency, auditability, and long-term reproducibility of analytical and machine learning outputs.

Every derived feature is traceable back to authoritative source datasets through deterministic transformation logic. No feature is introduced without an explicitly defined origin, transformation path, and intended analytical meaning.

D.1 Feature Lineage Principles

Feature lineage within this project is governed by the following principles:

- Every feature must be traceable to one or more authoritative source tables
- Transformation logic must be deterministic and reproducible
- Biological annotations and analytical abstractions must remain clearly separated
- Features must not embed implicit clinical or causal assumptions

Lineage tracking is enforced through explicit table dependencies, controlled joins, and versioned transformation logic rather than automated enterprise lineage tooling. This approach aligns with the project's local-first execution constraints while preserving scientific rigor.

D.2 Example Feature Lineage

This section illustrates a representative end-to-end feature lineage example.

Feature Name

`disease_gene_count`

Source Tables

- `silver.gene_disease_links`
- `silver.diseases`

Transformation Logic

Distinct gene identifiers are grouped by normalized disease identifier. The total number of unique genes associated with each disease is computed deterministically.

Biological Interpretation

This feature quantifies the genetic complexity of a disease by measuring how many distinct genes contribute to its etiology.

Downstream Usage

- Polygenic risk classification
- Disease complexity stratification
- Feature input for classical and graph-based machine learning models

D.3 Lineage Enforcement Strategy

Feature lineage enforcement is achieved through structural and procedural safeguards rather than post hoc validation.

- Consistent naming conventions distinguishing biological entities from analytical features
- Explicit separation between silver-layer biological tables and gold-layer analytical tables
- One-directional data flow preventing analytical features from influencing upstream biological representations
- Version-controlled transformation logic executed in fixed notebook or script sequences

This strategy prevents semantic drift, feature leakage, and untraceable transformations as the system evolves.

D.4 Lineage Scope and Limitations

Feature lineage documentation in this appendix focuses on analytical and machine learning readiness rather than regulatory-grade audit compliance. While the lineage framework is rigorous, it is intentionally lightweight and compatible with local execution environments.

Future extensions may incorporate automated lineage capture or metadata catalogs if the system transitions to production-grade infrastructure.

Appendix E

System Architecture Reference (Logical)

This appendix provides a consolidated logical view of the system architecture underpinning the DNA Gene Mapping Project. It serves as a stable architectural reference that complements the detailed, phase-wise descriptions provided in the main chapters.

The architecture described here is logical rather than deployment-specific. It reflects conceptual responsibilities, data flow boundaries, and interaction patterns without assuming cloud infrastructure or managed services.

E.1 Logical Architecture Layers

The system is organized into a set of clearly defined logical layers, each with a strict responsibility boundary. These layers collectively support data ingestion, enrichment, analytics, machine learning readiness, and user-facing consumption.

- **Data Ingestion Layer** Responsible for acquiring authoritative public biomedical datasets, validating file integrity, and performing minimal structural preprocessing. This layer ensures that upstream source data is captured faithfully without introducing analytical assumptions.
- **Normalization and Enrichment Layer** Transforms raw inputs into structured biological entities such as genes, variants, diseases, transcripts, and proteins. Identifier normalization, ontology alignment, and biological enrichment occur at this stage.
- **Feature Engineering Layer** Derives analytical and machine learning features from normalized entities. This includes aggregation, scoring, and classification logic designed for interpretability and traceability rather than model optimization.

- **Machine Learning Layer** Consumes curated feature tables to support classical machine learning, graph-based modeling, and deep learning feature generation. This layer is explicitly optional and modular, enabling selective experimentation without impacting upstream data integrity.
- **Visualization and Application Layer** Exposes analytical outputs through dashboards, reports, and programmatic APIs. This layer is strictly read-only with respect to curated biological and analytical data.

Each layer interacts only with its immediate neighbors, preventing tight coupling and uncontrolled dependency growth.

E.2 Technology Stack Summary

The technology stack is intentionally constrained to local and community-grade tooling to ensure feasibility, transparency, and reproducibility under realistic hardware limitations.

- **Databricks Community Edition** Used as the primary analytical execution environment for Spark-based transformations, feature engineering, and exploratory machine learning workflows. Databricks is treated as a computation workspace rather than a managed cloud platform.
- **Apache Spark (PySpark and Spark SQL)** Provides distributed processing capabilities, schema-aware transformations, and scalable joins required for large genomic datasets.
- **Delta-style Managed Tables** Used conceptually to enforce schema stability and versioned persistence within the Databricks environment, without reliance on external object storage.
- **Python Scientific Ecosystem** Supports orchestration, feature logic, and machine learning experimentation using well-established libraries.
- **SQL** Serves as the primary analytical access layer for downstream reporting, validation, and BI integration.
- **Local Analytical Tools** Power BI Desktop and optional local PostgreSQL are used strictly for consumption and exploration, not for primary data processing.

This stack prioritizes correctness, explainability, and auditability over deployment scale.

E.3 Architectural Scope and Boundaries

This architectural reference deliberately excludes:

- Managed cloud storage services
- Auto-scaling compute clusters
- Real-time streaming pipelines
- Production-grade CI/CD systems

These exclusions are intentional design decisions aligned with the project's research-first and local-execution philosophy. The logical architecture remains valid even if future deployment strategies evolve.

Appendix F

Governance and Compliance Reference

F.1 Regulatory Disclaimer

This system is intended solely for research, educational, and exploratory data analysis purposes. It is not designed, validated, or approved for clinical diagnosis, medical decision-making, patient treatment, or regulatory submission.

All analytical outputs, models, and visualizations generated by this system must not be interpreted as medical advice or used as a substitute for professional clinical judgment.

F.2 Ethical Safeguards

The project adheres to strict ethical safeguards to ensure responsible handling and interpretation of biomedical data.

- No patient-identifiable, personal, or protected health information is used at any stage of the system.
- All datasets originate from publicly available biomedical research resources released for non-clinical use.
- No inference of personal identity, ancestry, or protected attributes is performed.
- Analytical logic and feature derivations are designed to remain transparent and explainable.

Ethical considerations are enforced at the data modeling, feature engineering, and analytical layers to prevent misuse or over-interpretation of results.

F.3 Data Licensing and Attribution

All external data sources incorporated into the project are reviewed for licensing terms, attribution requirements, and redistribution constraints prior to ingestion.

Key principles followed include:

- Use of publicly available biomedical research datasets
- Respect for source-specific attribution and citation requirements
- No redistribution of restricted or proprietary raw data
- Preference for derived analytical features over raw source content where licensing constraints apply

Licensing considerations are treated as first-class governance requirements rather than post-processing concerns.

F.4 Governance Enforcement Strategy

Governance is enforced through architectural and process-level controls rather than policy statements alone.

- Read-only access to curated analytical and feature tables
- Explicit schema contracts across reference, silver, and gold layers
- Deterministic feature generation with documented lineage
- Clear separation between biological annotations and analytical interpretations

These controls ensure compliance, traceability, and ethical boundaries remain intact as the system evolves.

F.5 Scope Boundaries and Intended Use

The project explicitly avoids:

- Patient-level inference or individualized risk assessment
- Clinical decision automation
- Diagnostic or therapeutic recommendations
- Regulatory claims or validation assertions

The system is designed to demonstrate rigorous biomedical data engineering, analytics readiness, and research-grade modeling practices within clearly defined non-clinical boundaries.

Appendix G

System Architecture (Local Execution)

G.1 End-to-End System Architecture (Local)

This project is implemented entirely using local and community-grade tools without reliance on any commercial cloud infrastructure. All data processing, transformation, and analysis are performed within Databricks Community Edition operating in a local execution context.

Databricks is used strictly as an analytical processing environment and not as a managed cloud platform. It provides Spark-based computation, schema-aware table management, and notebook-driven orchestration without introducing cloud dependencies.

Raw genomic and clinical datasets are stored as flat files on the local filesystem and accessed directly by Spark jobs. Data persistence is handled through Databricks-managed tables within the local workspace. Optional integration with a local PostgreSQL instance is supported for analytical access and external tooling but is not required for core pipeline execution.

This architecture emphasizes reproducibility, transparency, and feasibility under constrained hardware environments, aligning with academic and independent research workflows.

G.2 Data Layer Architecture (Logical, Not Cloud-Based)

The Bronze–Silver–Gold data layering model is applied as a logical organizational strategy rather than a cloud storage pattern. No object storage systems, cloud data lakes, or managed Delta Lake services are used.

- The Bronze layer represents raw or transient ingestion data and may be omitted where unnecessary.
- The Silver layer contains cleaned, normalized, and biologically enriched entities suitable for deterministic analytics.
- The Gold layer stores feature-engineered datasets intended for analytics and machine learning consumption.

All layers exist as Spark tables within Databricks Community Edition and are backed by local storage mechanisms. The layering model improves data clarity, lineage tracking, and analytical stability without introducing cloud-specific assumptions.

G.3 Feature Lineage and Dependency Tracking

Feature lineage within the system is tracked through explicit table dependencies, deterministic transformation logic, and structured notebook execution order rather than enterprise lineage tooling.

Each Silver and Gold table is derived from well-defined upstream datasets using controlled joins and documented transformation steps. Lineage enforcement relies on:

- Versioned notebooks and scripts
- Explicit table naming conventions
- Stable primary keys and schema contracts
- Separation of biological normalization and analytical feature logic

This approach ensures traceability and auditability while remaining compatible with local execution environments.

G.4 Gene–Variant–Disease Relationship Model

The gene–variant–disease relationship model represents the biological associations captured within the system and is independent of deployment infrastructure.

Genomic variants are mapped to genes using curated reference annotations. Genes are associated with diseases through validated ontology-aligned mappings. Disease hierarchies are preserved to support analysis across multiple levels of clinical granularity.

This model is implemented using relational tables and Spark DataFrames without reliance on graph databases or cloud-based biomedical platforms. Network-based analyses

are supported through explicit relationship tables rather than infrastructure-dependent graph services.

The relationship model serves as the conceptual foundation for downstream analytics, feature engineering, and machine learning workflows.

Appendix H

Machine Learning Architecture (Local Research Mode)

This appendix defines the machine learning architecture used in the DNA Gene Mapping Project under a strictly local and research-oriented execution model. The architecture is designed to support analytical experimentation, feature generation, and model interpretability without reliance on cloud infrastructure, GPU acceleration, or managed machine learning services.

The machine learning layer is positioned as an extension of the curated biological data layers rather than as an isolated modeling system. All models consume schema-stable, versioned feature tables and produce reusable analytical features rather than opaque predictions.

H.1 Architectural Scope and Design Principles

The machine learning architecture adheres to the following core principles:

- Models consume curated Silver and Gold layer tables only
- No model operates directly on raw or partially normalized data
- Outputs are persisted as features, not final decisions
- Interpretability and reproducibility are prioritized over scale
- All execution is compatible with CPU-only environments

Machine learning is treated as a controlled analytical enhancement layer that augments biological insight while preserving full traceability to upstream data sources.

H.2 End-to-End Machine Learning Flow

Classical machine learning, graph-based learning, and deep learning models operate independently but follow the same architectural contract: consume validated features and emit model-derived representations suitable for reintegration into Gold-layer analytical tables.

Figure H.1 illustrates the logical flow of machine learning within the local research environment. Curated biological entities from the Reference and Silver layers are transformed into analytical feature tables in the Gold layer. These feature tables serve as the sole input to downstream modeling workflows.

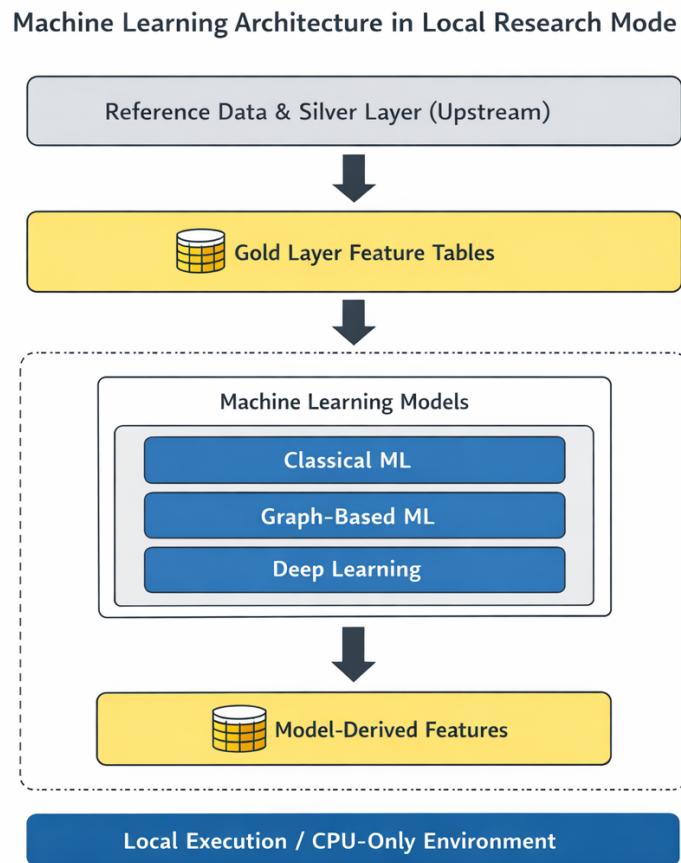


Figure H.1: Logical machine learning architecture in local research mode. Curated Silver and Gold layer features feed classical, graph-based, and deep learning models. All models operate in a CPU-only environment and persist outputs as reusable analytical features rather than final predictions.

H.3 Model Categories Supported

The architecture supports three complementary categories of machine learning models:

- Classical Machine Learning** Tabular models such as logistic regression, tree-based

methods, and ranking models are used for interpretable tasks including gene prioritization, disease complexity classification, and variant risk scoring.

Graph-Based Machine Learning Graph-oriented models operate on curated gene–variant–disease relationships to capture network effects and polygenic structure that are not visible in flat tables.

Deep Learning Models Neural architectures such as CNNs, GCNs, and transformers are introduced selectively as feature generators. These models learn representations from sequences, networks, or ordered biological contexts while remaining optional and non-blocking to the core pipeline.

H.4 Local Execution Constraints

All machine learning workflows are explicitly designed to operate under the following constraints:

- CPU-only execution
- No distributed training clusters
- Limited memory and storage availability
- No external APIs or managed services

These constraints are treated as deliberate architectural boundaries rather than limitations. They enforce disciplined feature engineering, controlled model complexity, and reproducible experimentation aligned with real-world research environments.

H.5 Feature Persistence and Governance

Model outputs are persisted exclusively as structured features within Gold-layer tables. Raw predictions, intermediate tensors, and transient artifacts are not treated as durable outputs.

Each model-derived feature maintains:

- Explicit linkage to the generating model type
- Versioned feature definitions
- Clear biological or analytical interpretation

This approach ensures that machine learning enhances analytical depth without compromising auditability, explainability, or long-term maintainability of the system.

Appendix I

Analytics and Visualization Layer (Local)

This appendix defines the analytics and visualization layer of the DNA Gene Mapping Project as implemented under a strictly local execution model. The purpose of this layer is to expose curated analytical insights derived from Gold-layer feature tables in a form that is interpretable, auditable, and suitable for exploratory and decision-support analysis.

The visualization layer is explicitly decoupled from data ingestion, transformation, and machine learning workflows. It operates only on validated, schema-stable analytical outputs and does not perform any data mutation or feature derivation.

I.1 Analytical Consumption Strategy

All analytical consumption is performed on Gold-layer tables produced by deterministic feature engineering pipelines. These tables are designed to be directly consumable by analytics tools without requiring additional joins, enrichment logic, or preprocessing.

Key principles governing analytical consumption include:

- Read-only access to curated analytical tables
- No direct interaction with Silver or Reference layers
- Explicit separation between analytical computation and visualization
- Full traceability from visualized metrics back to feature definitions

This strategy ensures that all visual outputs reflect validated analytical state rather than intermediate or transient results.

I.2 Power BI Desktop Architecture

Power BI is used exclusively in Desktop mode as a local visualization and exploratory analytics tool. No Power BI Service, cloud publishing, gateways, or online sharing mechanisms are used at any stage.

Data access is performed using:

- Exported analytical tables (CSV or Parquet)
- Direct local database connections where applicable

Power BI dashboards are designed to consume denormalized Gold-layer tables, ensuring predictable performance and eliminating dependency on live transformation logic.

I.3 Dashboard Design Principles

Dashboards are structured around analytical questions rather than raw datasets. Visual components are grouped by biological and analytical themes to support intuitive exploration by non-engineering users.

Core design principles include:

- Disease-centric views emphasizing complexity and burden
- Gene-centric views supporting prioritization and comparison
- Variant-centric views highlighting pathogenicity, rarity, and confidence
- Consistent filtering across disease, gene, and variant dimensions

All metrics displayed in dashboards correspond directly to documented Gold-layer features and retain semantic consistency across views.

I.4 Data Refresh and Governance

The analytics layer does not perform live data refresh or background synchronization. Data updates occur only when upstream analytical pipelines are re-executed and validated.

Governance rules enforced at the visualization layer include:

- No manual data modification within dashboards
- No calculated fields that replicate feature engineering logic
- Explicit version awareness of analytical datasets

This approach ensures that visual analytics remain reproducible, auditable, and aligned with documented feature definitions.

I.5 Local-Only Execution Constraints

The analytics and visualization layer is intentionally designed to operate fully offline using local compute and storage resources. This constraint enforces disciplined feature design and prevents reliance on external services or managed visualization platforms.

Local execution ensures:

- Zero operational cost
- Full control over sensitive biomedical analytics
- Reproducibility independent of external infrastructure

These constraints align the visualization layer with the overall research-first, local-execution philosophy of the project.

Appendix J

Deep Learning Architecture (Conceptual)

This appendix defines the conceptual deep learning architecture of the DNA Gene Mapping Project. It describes how neural models conceptually integrate with curated biological data layers and contribute derived representations to downstream analytics without prescribing concrete execution pipelines or training schedules.

The architecture documented here is intentionally abstract, model-agnostic, and research-oriented. It establishes deep learning readiness while preserving interpretability, feature lineage, and compatibility with constrained local execution environments.

J.1 Architectural Scope and Design Principles

The deep learning layer is positioned as an analytical extension of the Silver and Gold data layers. Neural models do not operate on raw biological data; instead, they consume curated, schema-stable tables to ensure correctness, traceability, and reproducibility.

The conceptual architecture follows these principles:

- Neural models consume only curated Silver- or Gold-layer tables
- Model outputs are persisted as derived features, not final predictions
- Deep learning augments deterministic biological enrichment rather than replacing it
- No model is required to be end-to-end or continuously trained
- Interpretability and feature lineage are prioritized over model complexity

J.2 Conceptual Deep Learning Flow

Figure J.1 illustrates the high-level conceptual flow of deep learning within the system. Curated biological entities are transformed into structured numerical representations suitable for neural modeling. Multiple model families may operate independently on the same curated inputs.

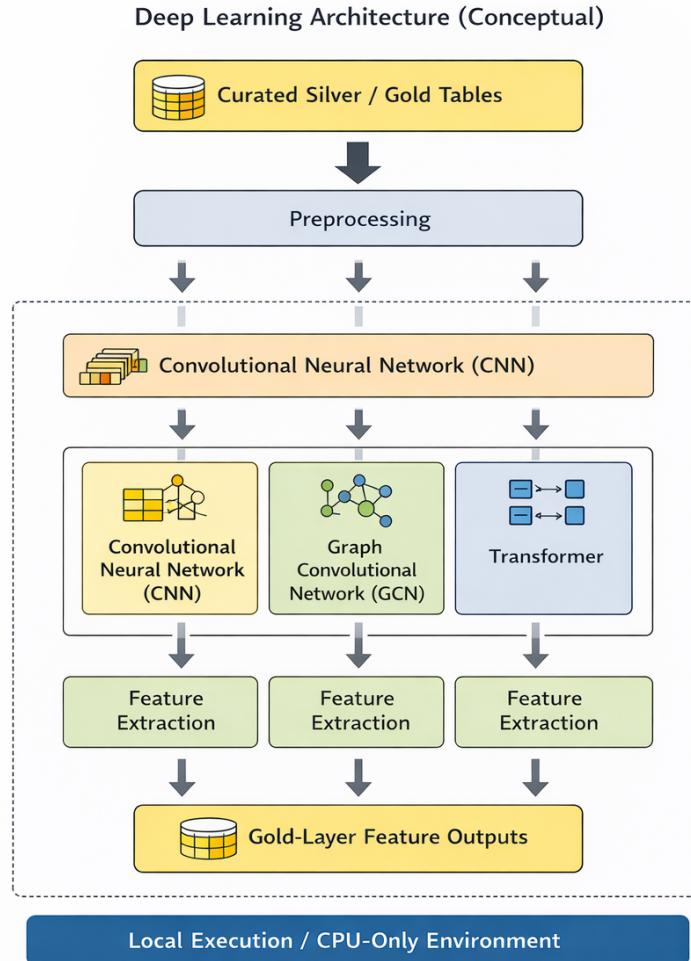


Figure J.1: Conceptual deep learning architecture showing how curated Silver and Gold tables feed CNN, GCN, and Transformer models. Neural outputs are extracted as structured features and persisted into Gold-layer feature tables under local, CPU-only execution constraints.

J.3 Neural Model Families and Conceptual Roles

Convolutional Neural Networks are conceptually used to capture localized biological patterns such as conserved sequence motifs, protein functional regions, and clustered variant effects that are difficult to encode using rule-based logic.

Graph Convolutional Networks are conceptually applied to curated biological graphs, enabling the propagation of information across gene–protein–disease relationships and

supporting network-aware feature generation.

Transformer architectures are conceptually included to model long-range dependencies across ordered biological data such as protein sequences, variant ordering within genes, or longitudinal biological signals.

J.4 Feature-Oriented Integration Philosophy

All deep learning models contribute intermediate representations rather than final decisions. These representations are transformed into structured analytical features and persisted into Gold-layer tables.

This approach ensures that:

- Neural outputs remain reusable across multiple analytical workflows
- Feature lineage remains explicit and auditable
- Classical and neural features can be combined transparently

J.5 Local Execution Constraints

The conceptual architecture is explicitly designed for local execution:

- CPU-only processing
- No distributed training infrastructure
- Small batch sizes and shallow configurations
- No reliance on external managed services

These constraints are deliberate design choices aligned with research-first, cost-aware, and reproducible system objectives.

Appendix K

Deep Learning Execution Architecture

This appendix documents the execution-level architecture for deep learning within the DNA Gene Mapping Project. It provides structural reference pipelines for how convolutional, graph-based, and transformer models are executed under local constraints and how their outputs are integrated into analytical feature tables.

The architectures described here are optional extensions to the core pipeline and are not required for system correctness or completeness.

K.1 Execution Overview

Figure K.1 presents the execution-level flow for deep learning. Curated Silver and Gold tables are transformed into numerical inputs, processed by neural models, and converted into interpretable Gold-layer features.

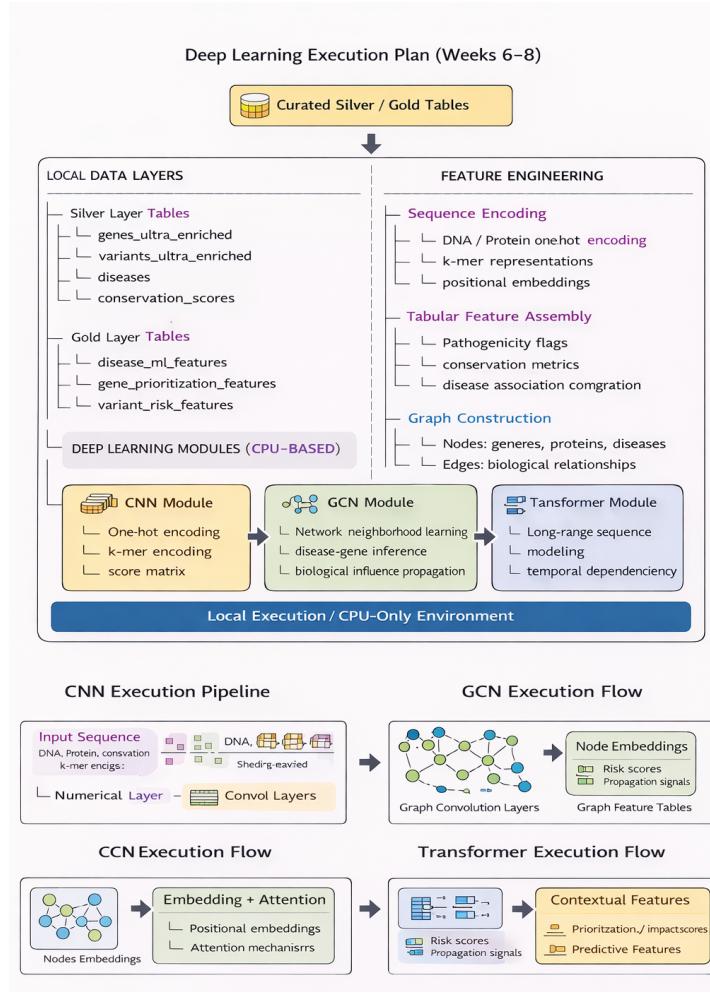
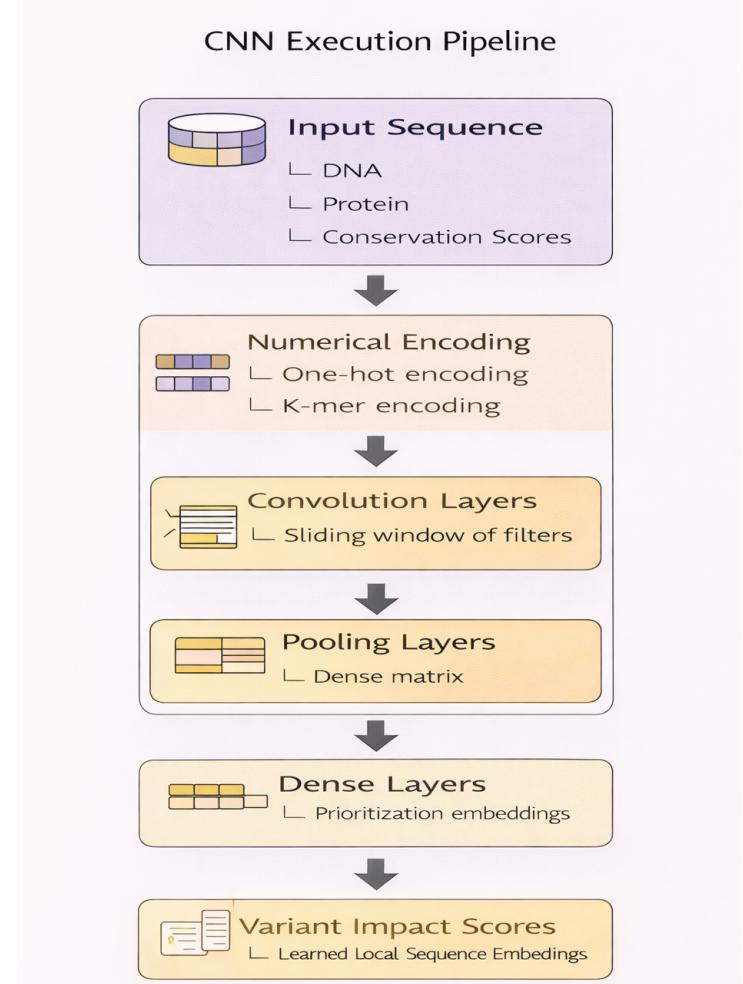


Figure K.1: Deep learning execution flow showing curated Silver and Gold tables feeding CNN, GCN, and Transformer architectures, with learned representations persisted as interpretable Gold-layer features under local execution constraints.

K.2 CNN Execution Pipeline

Convolutional Neural Networks are executed on fixed-size numerical inputs derived from curated biological sequences and conservation matrices.



CPU-based deep learning execution plan illustrating local STE and Gold data layers, feature engineering steps, CNN, GCN, and Transformer models, and output persistence into interpretable Gold-layer feature tables.

Figure K.2: CNN execution pipeline illustrating sequence encoding, convolutional feature extraction, pooling, and dense layers used to derive interpretable variant- and protein-level feature representations.

CNN-derived outputs include:

- Local sequence motif scores
- Structural pattern indicators
- Conservation-aware embeddings

K.3 GCN Execution Pipeline

Graph Convolutional Networks operate on explicitly constructed biological graphs derived from curated relational tables.

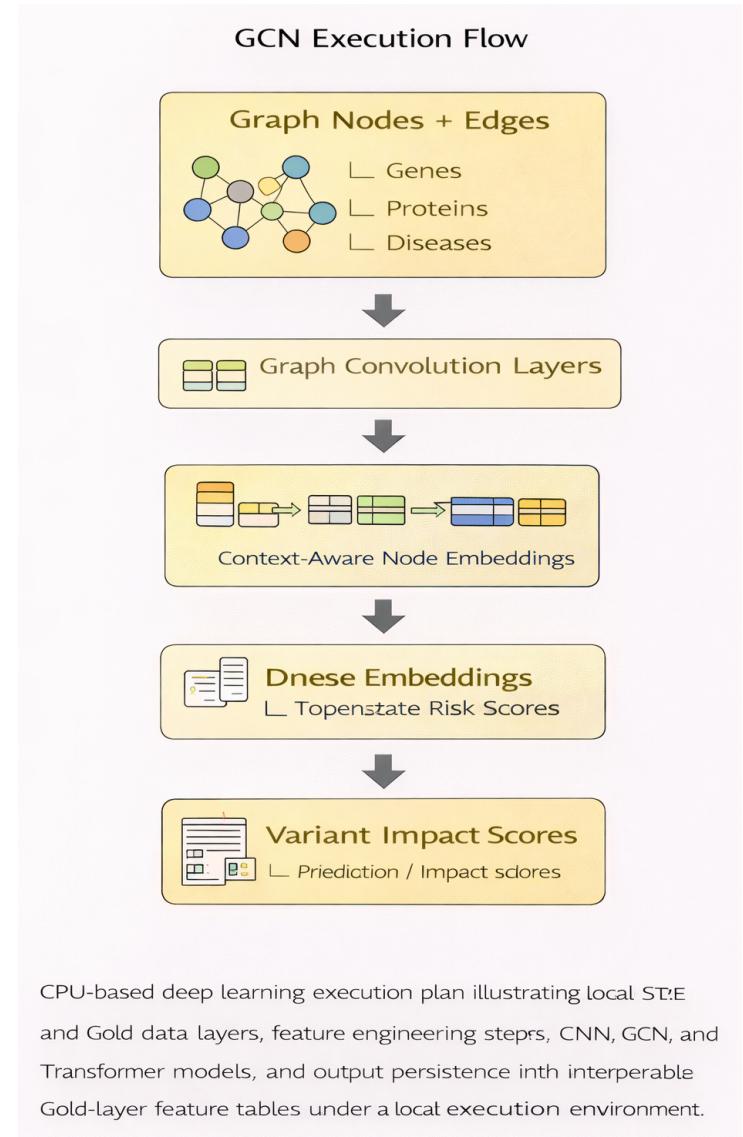


Figure K.3: GCN execution flow showing curated graph construction, neighborhood aggregation, node embedding generation, and disease–gene risk signal derivation.

GCN-derived features include:

- Gene network importance scores
- Disease association propagation metrics
- Polygenic influence indicators

K.4 Transformer Execution Pipeline

Transformer models are executed on tokenized and ordered biological sequences derived from curated inputs.

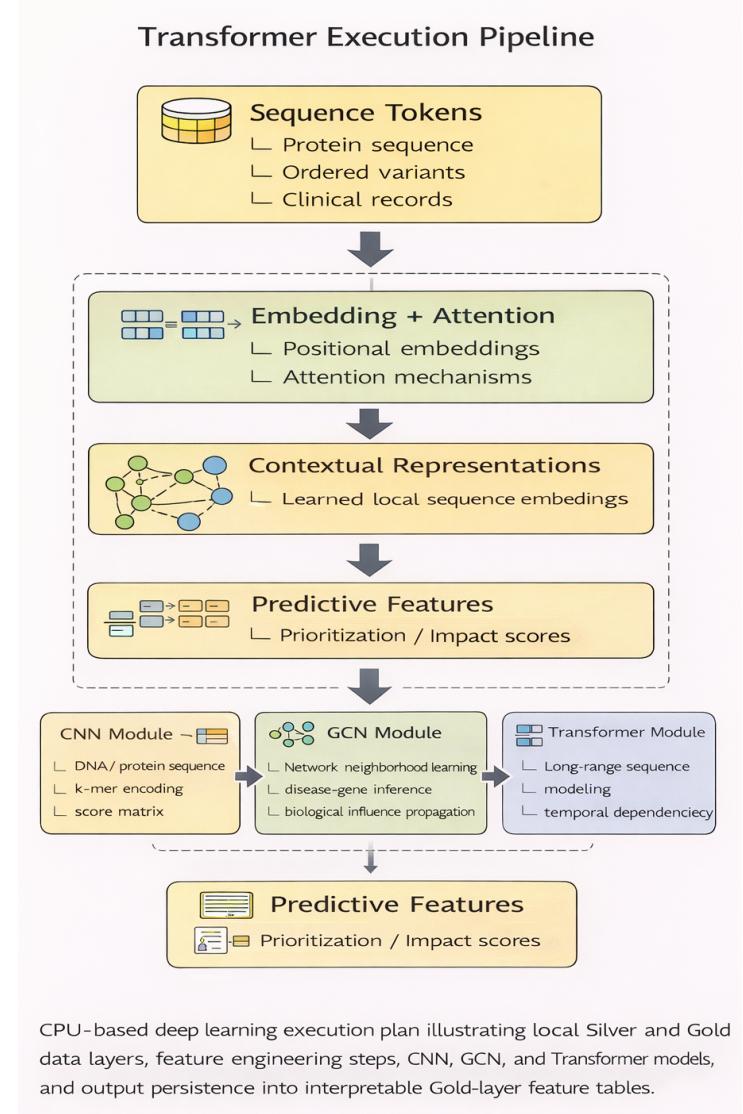


Figure K.4: Transformer execution pipeline showing tokenization, attention-based contextual embedding, and feature extraction for downstream analytical use.

Transformer-derived outputs include:

- Context-aware sequence features
- Long-range dependency indicators
- Temporal or positional risk signals

K.5 Feature Integration Rules

All deep learning outputs are persisted as derived Gold-layer features and never overwrite curated biological annotations.

Integration rules include:

- Explicit versioning of learned features
- Traceable linkage to source data and model architecture
- Separation of biological facts from learned representations
- Compatibility with classical machine learning pipelines

K.6 Execution Constraints

All execution pipelines described here operate under:

- CPU-only processing
- Limited memory availability
- Small batch sizes
- No external managed ML platforms

These constraints preserve transparency, reproducibility, and feasibility within local research environments.

Appendix L

Model-to-Feature Mapping Inventory

This appendix documents how outputs generated by machine learning and deep learning models are transformed into structured, reusable analytical features within the DNA Gene Mapping Project. The system does not treat raw model predictions as final outputs. Instead, learned representations and derived signals are persisted as interpretable features within the Gold layer.

The purpose of this appendix is to ensure transparency, auditability, and long-term reuse of model-derived features across analytical, research, and exploratory workflows.

L.1 Purpose and Mapping Principles

All model-to-feature mappings in this project follow a consistent set of principles:

- Model outputs are never stored as opaque predictions
- Learned representations are converted into named, versioned features
- Features preserve biological interpretability wherever possible
- Each feature is traceable to a specific model type and input data source
- No model-derived feature overwrites curated biological annotations

Model-derived features are treated as analytical signals rather than biological ground truth.

L.2 CNN-Derived Feature Mapping

Convolutional Neural Networks are used to learn localized patterns from biological sequences and structured matrices. Intermediate representations and aggregation outputs are converted into variant- and protein-level features.

Table L.1: CNN model outputs mapped to analytical features

CNN Type	Output	Mapped Feature Name	Target Gold Table
Intermediate convolution activations	sequence motif score	gold.variant_risk_features	
Pooling layer embeddings	variant_local_pattern_strength	gold.variant_risk_features	
Protein-region embeddings	protein_region_embedding	gold.protein_features	
Conservation-aware filters	cnn_derived_conservation_indicator	gold.variant_risk_features	

These features capture localized biological patterns while remaining suitable for downstream statistical analysis and classical machine learning models.

L.3 GCN-Derived Feature Mapping

Graph Convolutional Networks operate on curated biological graphs composed of genes, proteins, variants, and diseases. Learned node embeddings and propagation scores are converted into network-aware analytical features.

Table L.2: GCN model outputs mapped to analytical features

GCN Type	Output	Mapped Feature Name	Target Gold Table
Gene node embeddings	gene_node_embeddings	gene_network_importance	gold.gene_prioritization_features
Disease node embeddings	disease_node_embeddings	disease_association_strength	gold.disease_ml_features
Edge influence scores	edge_influence_scores	polygenic_influence_score	gold.disease_ml_features
Neighborhood aggregation metrics	neighborhood_aggregation_metrics	network_propagation_signal	gold.gene_prioritization_features

These features encode relational and polygenic effects that are not observable through tabular analysis alone.

L.4 Transformer-Derived Feature Mapping

Transformer architectures are used to model long-range dependencies within ordered biological sequences and contextual feature sets. Attention-weighted representations are transformed into contextual analytical features.

Table L.3: Transformer model outputs mapped to analytical features

Transformer Output Type		Mapped Feature Name	Target Gold Table
Contextual embeddings	sequence_context_embeddings	sequence_context_score	gold.variant_risk_features
Attention-weighted representations	attention_weighted_representations	long_range_dependency_index	gold.variant_risk_features
Ordered sequence modeling outputs	ordered_sequence_outputs	temporal_risk_signal	gold.disease_ml_features
Protein sequence context vectors	protein_sequence_context_vectors	protein_contextual_embedding	gold.protein_features

Transformer-derived features provide global context signals that complement localized and network-based representations.

L.5 Feature Governance Rules

All model-derived features must comply with the following governance rules:

- Features are versioned and reproducible
- Source model architecture and inputs are documented
- Feature definitions remain stable across pipeline runs
- Deprecated features are retained for auditability
- Model-derived features never replace curated annotations

This governance framework ensures that analytical insights remain interpretable, defensible, and reusable across future iterations of the system.

Appendix M

Training, Evaluation, and Validation Strategy

This appendix defines the principles, methodologies, and safeguards used for training, evaluating, and validating machine learning and deep learning models within the DNA Gene Mapping Project. The strategies described here prioritize biological plausibility, interpretability, and reproducibility over raw predictive performance.

All model development is conducted under constrained local execution conditions and is designed for exploratory research and analytical feature generation rather than clinical deployment.

M.1 Training Strategy

Model training follows a conservative and controlled approach designed to minimize overfitting, leakage, and biological misinterpretation.

Key training principles include:

- Use of curated Silver and Gold layer tables only
- Explicit separation between identifiers, labels, and features
- Stratified sampling where class imbalance exists
- Limited training epochs to prevent memorization
- Reproducible initialization and fixed random seeds

Training workflows favor incremental experimentation, allowing models to be evaluated and refined without destabilizing downstream analytical features.

M.2 Evaluation Metrics

Evaluation emphasizes clinically and biologically meaningful metrics rather than purely statistical optimization.

Variant-Level Models

Evaluation metrics include:

- Precision and recall for pathogenic and benign classifications
- Class stability across disease groups
- Separation between pathogenic and uncertain variants

Gene-Level Models

Evaluation metrics include:

- Ranking consistency for clinically relevant genes
- Sensitivity to pathogenic variant burden
- Robustness across monogenic and polygenic diseases

Disease-Level Models

Evaluation metrics include:

- Accuracy of disease complexity classification
- Stability across disease categories
- Interpretability of risk and burden indicators

M.3 Validation Approach

Validation is performed using a combination of quantitative checks and biological reasoning rather than automated acceptance thresholds.

Validation strategies include:

- Cross-validation with controlled splits to prevent gene or disease leakage
- Comparison against curated clinical annotations
- Manual inspection of high-impact features and outliers

- Consistency checks across multiple model runs

No automated clinical decisions are derived from model outputs. All results are treated as analytical signals requiring contextual interpretation.

M.4 Risk and Limitations

Several limitations are explicitly acknowledged:

- Limited training data for rare diseases
- Bias introduced by reporting and curation practices in public databases
- Absence of GPU acceleration
- Small batch sizes and constrained memory

These limitations are mitigated through conservative modeling choices, feature transparency, and explicit documentation rather than algorithmic complexity.

M.5 Final Integration Rule

A model is considered suitable for integration into the analytical pipeline only if:

- It produces stable and reproducible features
- Feature outputs are biologically interpretable
- Downstream analytical value is demonstrable
- Results remain consistent across repeated executions

Models that fail to meet these criteria are excluded without affecting the stability of the core data pipeline.

Appendix N

Local Execution Constraints and Design Decisions

This appendix documents the hardware, software, and operational constraints under which the DNA Gene Mapping Project is intentionally designed and executed. These constraints are not incidental limitations but deliberate design conditions that shape architectural choices, data modeling strategies, and analytical workflows throughout the system.

The purpose of this appendix is to provide transparency and justification for decisions that prioritize correctness, reproducibility, and feasibility over scale or performance optimization.

N.1 Hardware Constraints

All development and execution are performed on commodity, laptop-class hardware. The system is explicitly designed to function within the following constraints:

- CPU-only execution (no GPU acceleration)
- Limited memory availability (approximately 8 GB RAM)
- Finite local storage capacity
- Single-node execution environment

These constraints preclude large-scale distributed training or real-time inference but enable full transparency and control over data processing and modeling logic.

N.2 Software Constraints

The project uses only community-grade and locally executable tools. No managed cloud services or proprietary platforms are required.

Key software constraints include:

- Databricks Community Edition for distributed data processing
- Apache Spark for transformation and feature engineering
- Local filesystem-based data storage
- Optional local PostgreSQL for analytical querying
- Power BI Desktop for visualization

Cluster scaling, background jobs, and managed machine learning services are intentionally excluded.

N.3 Data Volume Management Strategy

Upstream biomedical datasets can be large and heterogeneous. To remain compatible with local execution, the system employs selective data reduction strategies:

- Retention of biologically meaningful subsets only
- Filtering of low-information or redundant records
- Avoidance of raw archive persistence once validation is complete
- Incremental enrichment rather than monolithic joins

These strategies reduce memory pressure while preserving analytical integrity.

N.4 Compute-Aware Feature Engineering

Feature engineering workflows are designed to minimize unnecessary computation and memory usage.

Key design choices include:

- Explicit staging of transformations
- Avoidance of repeated full-table scans
- Controlled joins with stable primary keys

- Pre-aggregation of features where possible

Complex transformations are decomposed into deterministic steps to ensure both performance stability and reproducibility.

N.5 Deferred Infrastructure Components

Several enterprise-grade infrastructure components are intentionally deferred, including:

- Scalable API deployments
- Persistent production databases
- Automated CI/CD pipelines
- Real-time model serving

Deferral of these components allows the project to focus on data correctness, feature richness, and scientific validity without introducing unnecessary operational complexity.

N.6 Research-First Architecture Justification

The system is structured to support exploration, hypothesis testing, and iterative refinement rather than production deployment. This research-first orientation aligns with common workflows in academic research, early-stage biotech, and independent clinical analytics.

The local-first architecture enables:

- Full visibility into all data transformations
- Zero operational cost
- Reproducible experimentation
- Freedom from vendor lock-in

This approach ensures that architectural decisions remain transparent, defensible, and adaptable as future resources and requirements evolve.

Appendix O

Project Folder Structure and Codebase Organization

This appendix documents the canonical repository structure used to organize data engineering, analytics, machine learning, backend services, frontend applications, and deployment artifacts within a single unified codebase. The structure is designed to support clear separation of concerns while maintaining tight integration between schemas, features, and analytical outputs.

The repository follows a monorepo design to ensure consistency, traceability, and coordinated evolution across all system components.

O.1 Top-Level Repository Structure

Figure O.1 illustrates the top-level directory layout of the project repository.

The top-level directories include:

- `data/` for raw, processed, analytical, and ML-ready datasets
- `databricks_notebooks/` for Spark-based transformation workflows
- `ml/` for machine learning and deep learning pipelines
- `backend/` for API services
- `frontend/` for user-facing applications
- `deployment/` for containerization and environment configuration
- `documentation/` for project documentation and references

Supporting files such as configuration, environment variables, dependency definitions, and licensing information are maintained at the repository root.

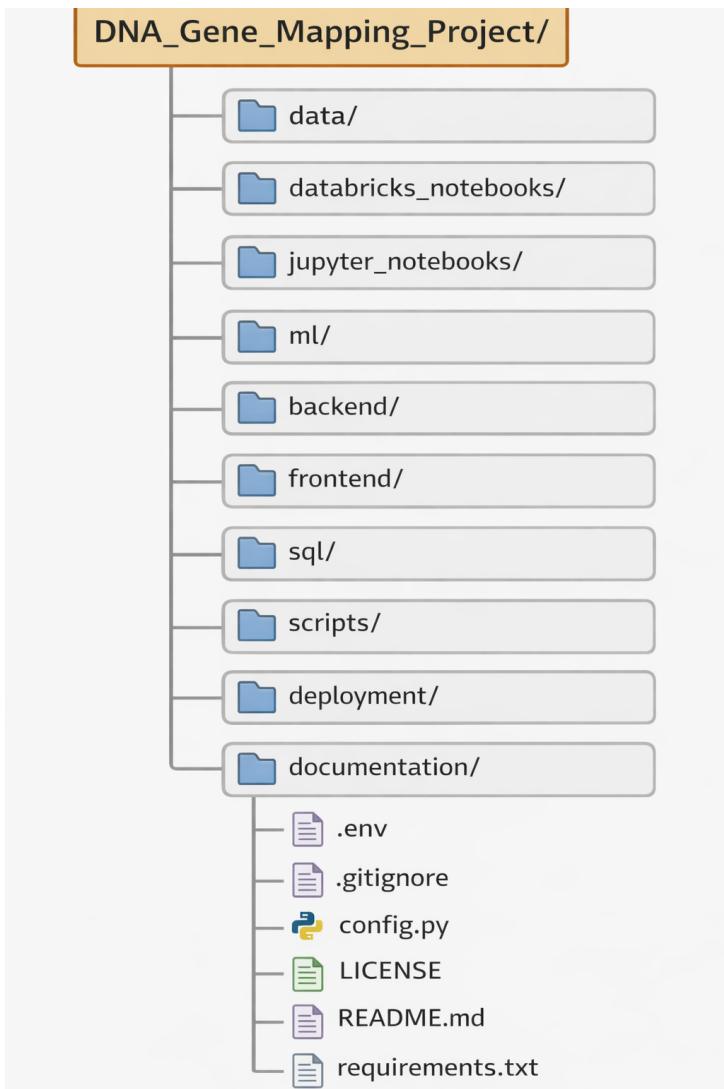


Figure O.1: Top-level monorepo structure showing separation between data, machine learning, backend services, frontend applications, and deployment artifacts.

O.2 Machine Learning and Deep Learning Workspace

The machine learning workspace encapsulates all feature engineering, modeling, training, and evaluation logic. Its internal structure is shown in Figure O.2.

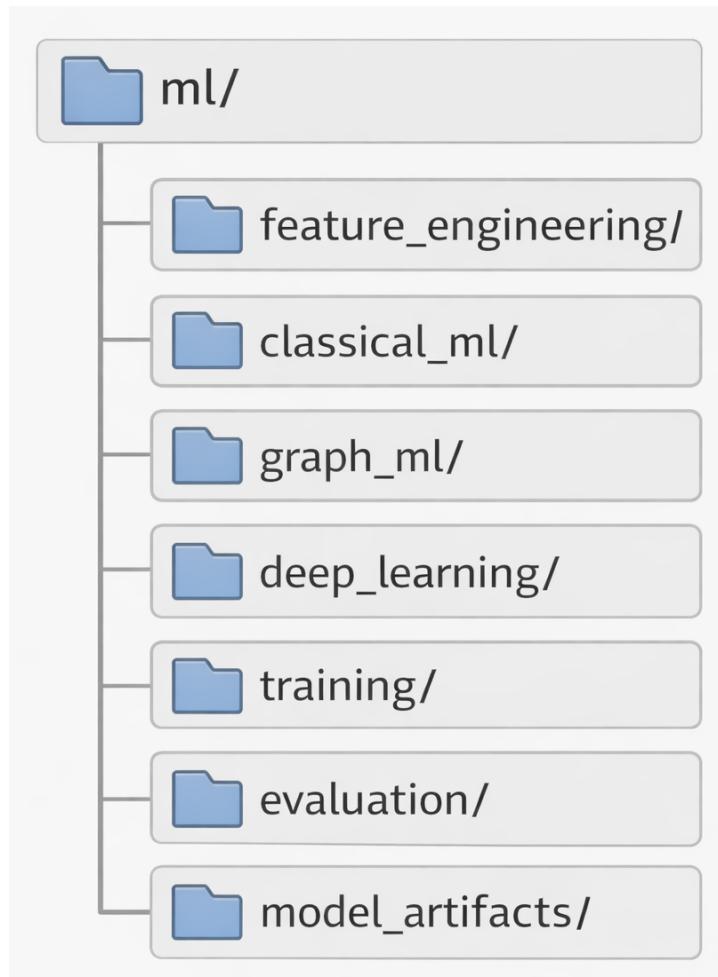


Figure O.2: Machine learning and deep learning workspace structure separating feature engineering, classical ML, graph-based models, deep learning, training, and evaluation.

This structure enables controlled experimentation while maintaining a clear boundary between feature generation, model training, and persisted artifacts.

O.3 Backend Application Structure

The backend application is implemented as a modular FastAPI service responsible for exposing read-only analytical and predictive endpoints. The internal organization is shown in Figure O.3.

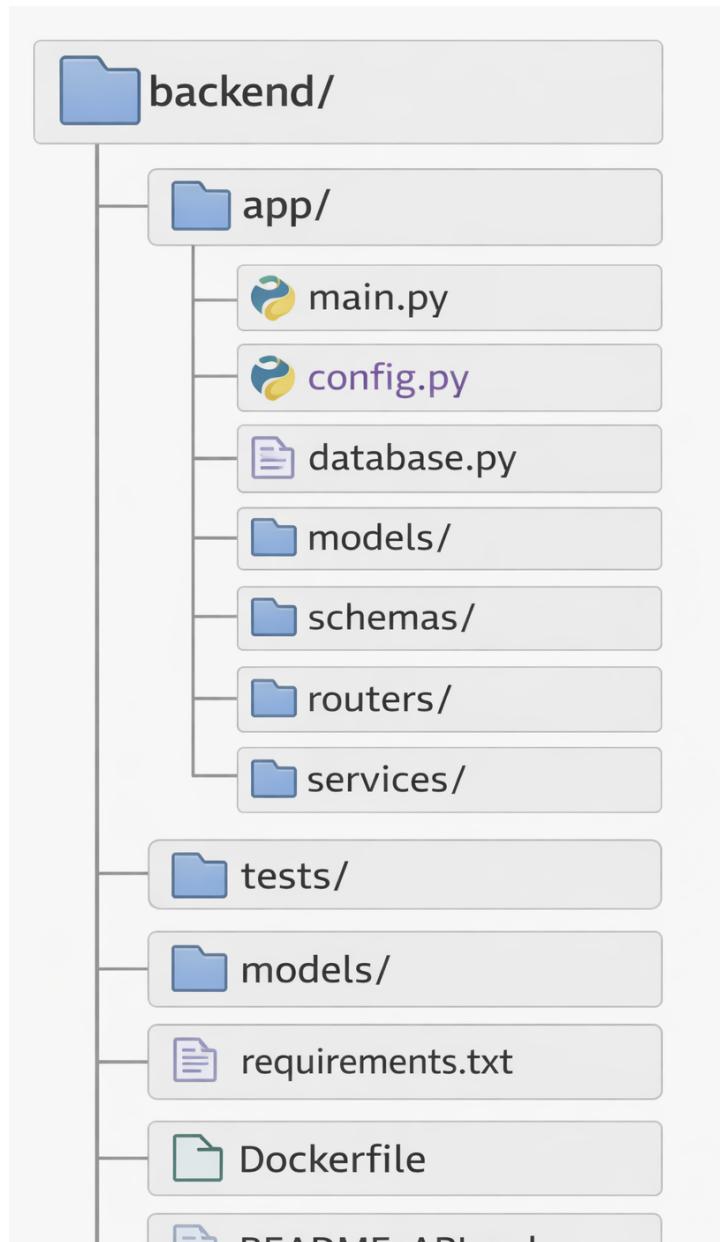


Figure O.3: Backend service structure illustrating application modules, routing, schema enforcement, service layers, and deployment artifacts.

Key design characteristics include:

- Explicit schema definitions for API contracts
- Separation between routing, services, and data access logic
- Read-only access to curated analytical data
- Containerized deployment support

O.4 Frontend Application Structure

The frontend is implemented as a single React-based single-page application composed of multiple analytical modules. Its structure is illustrated in Figure O.4.

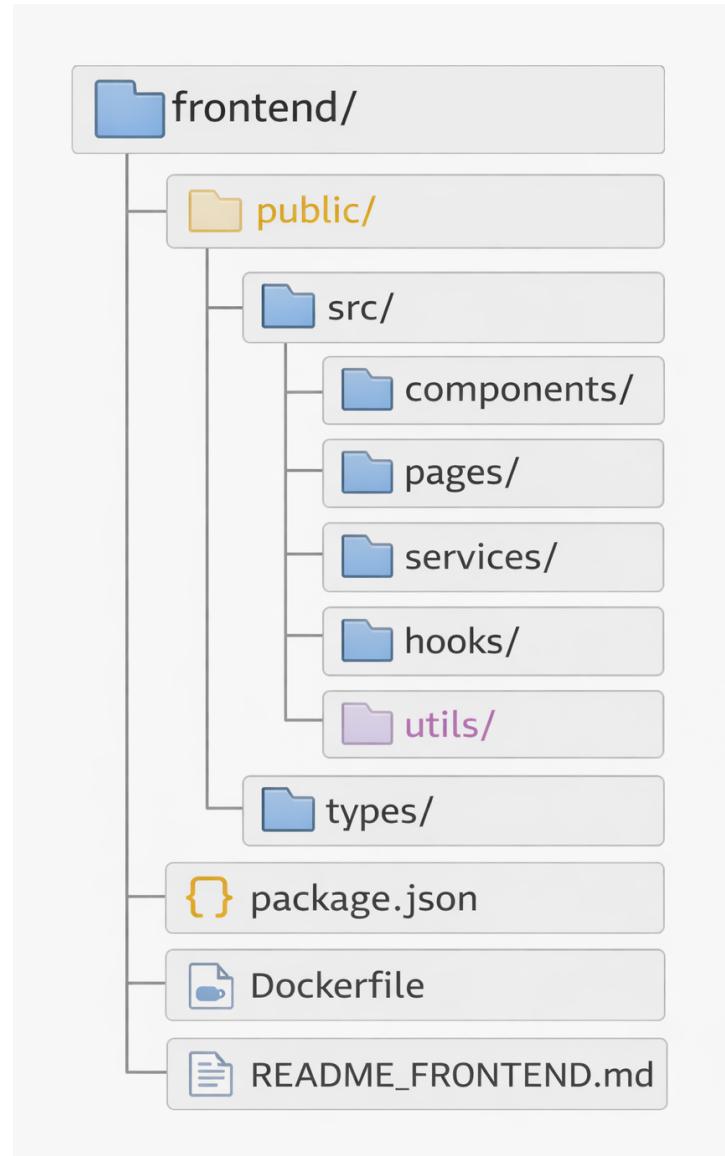


Figure O.4: Frontend application structure showing modular components, services, hooks, type definitions, and build configuration.

The frontend architecture emphasizes:

- Modular, reusable UI components
- Strong typing and contract alignment with backend APIs
- Separation between presentation logic and data access

O.5 Deployment and Environment Separation

Deployment-related artifacts are isolated from analytical and modeling code. The structure is shown in Figure O.5.

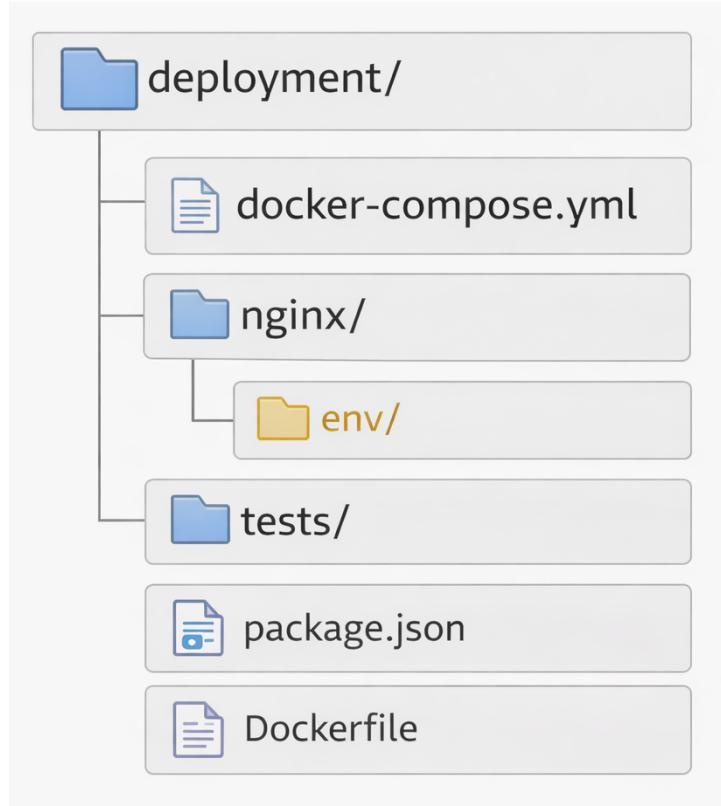


Figure O.5: Deployment directory structure containing container orchestration, reverse proxy configuration, and environment templates.

This separation ensures that development, staging, and future production deployments can be managed without modifying core data processing or modeling logic.

The overall repository organization supports maintainability, reproducibility, and long-term extensibility while remaining compatible with local-first execution constraints.