

BRYAN L.A.

LITERATURE NOTES BIOINFORMATICS

UNIVERSITY OF AMSTERDAM ∩ VRIJE UNIVERSITEIT

No copyright © 2025 Bryan L.A.

UNIVERSITY OF AMSTERDAM ∩ VRIJE UNIVERSITEIT

THIS FILE CONTAINS NOTES FROM THE LITERATURE PAPERS I READ THROUGHOUT MY BIOINFORMATICS
INTERNSHIP, 2025

Last updated June 2025

Contents

<i>June 11th /2025</i>	5
<i>June 12th /2025</i>	13
<i>June 13th /2025</i>	21
<i>June 15th /2025</i>	27
<i>June 16th /2025</i>	29
<i>June 19th /2025</i>	35
<i>June 20th /2025</i>	41
<i>June 25th /2025</i>	43
<i>June 27th /2025</i>	45
<i>June 28th /2025</i>	51
<i>Bibliography</i>	59

June 11th /2025

EVALUATION OF AN AUTOMATED GENOME INTERPRETATION MODEL FOR RARE DISEASE ROUTINELY USED IN A CLINICAL GENETIC LABORATORY ¹

Variant prioritization: Process of filtering and ranking a large number of genetic variants identified from sequencing data (i.e., exome/genome) to produce manageable shortlist of plausible candidates that may be responsible for a specific disease for a specific disease of phenotype.

Emedgene aims to reduce bottleneck by automatically generating a shortlist of candidate variants. Emedgene was evaluated on $n = 180$ retrospective *accuracy* previously solved exome cases and $n = 334$ prospective production cohort of consecutive clinical cases.

Correlated features with higher rank: Rare familial segregation, known pathogenicity, functional severity.

Accuracy was reduced in some cases due to incomplete genetic data (uncalled copy number variants) or atypical patient phenotypes. The AI-augmented analysis once integrated into workflow, achieved diagnostic rate 28.7% vs. comparable historical manual rates but significantly reduced the overall time required for case analysis by enabling a single cycle of review by a geneticist instead of two

Methods:

- Supervised learning approach, trained on dataset of 10^3 's of variants that had been manually curated
- Decision tree clustering. It creates model that - based on input - produces score for ranking variants
- **Features:** Integration of information from multiple sources
 - **Variant Level:** Allele freq/count and count of homozygotes in public (e.g., gnomAD) and internal databases
 - **Gene Level:** info about affected gene

¹ L Meng, R Attali, T Talmy, Y Regev, N Mizrahi, P Smirin-Yosef, L Vossaert, C Taborda, M Santana, I Machol, R Xiao, H Dai, C Eng, F Xia, and S Tzur. Evaluation of an automated genome interpretation model for rare disease routinely used in a clinical genetic laboratory. *Genet Med*, 25(6):100830, 2023. DOI: 10.1016/j.gim.2023.100830

- **Phenotypic Similarity:** Measure of match between patient's reported phenotypes (using Human Phenotype Ontology terms) and phenotypes associated with diseases linked to the variant's gene
- **family segregation:** Analysis of inheritance patterns (e.g., identifying *de novo* variants or assessing zygosity in context of recessive patterns in x family)
- **Functional Effect:** Predicted impact of variant on protein (e.g., loss-of-function effects like frameshift or nonsense mutations)
- **Known Pathogenicity:** Variants previously reported as pathogenic or likely pathogenic in databases like ClinVar or internal laboratory databases are given a very high rank

Training and evaluation:

- **Training:** Trained on manually curated variants to learn the correlations between the input features and the likelihood of a variant being diagnostic
- **Validation:** Model's performance was tested using CV on different segments of the data
- **Performance Metrics:** Sensitivity and specificity evaluate final model on accuracy cohort. sensitivity of 95.3% and a specificity of 99.9% for identifying a variant as a "most likely" candidate
- **Ranking Accuracy:** Rank of true diagnostic variant in the model's prioritized list of candidates

Current version does not account for certain types of genetic variation including CNVs, STRs, mitochondrial DNA variants

DEEP LEARNING-BASED RANKING METHOD FOR SUBGROUP AND PREDICTIVE BIOMARKER IDENTIFICATION IN PATIENTS²

Biomarkers associated with treatment effect heterogeneity = Predictive biomarkers. ML + Causal inference for predictive biomarker identification and ITR exploration.

To consider: Meta-learning, Q-learning, D-learning, DNNs for handling complex biomarker-treatment response relationship.

DeepRAB

mathematical framework to model *treatment effect heterogeneity* and construct *individualized Treatment Rule (ITR)*. Formulated as **supervised** ML problem where objective is to predict how much benefit

² Zhen Liu, Yifan Gu, and Xiaoyang Huang. Deep learning-based ranking method for subgroup and predictive biomarker identification in patients. *Communications Medicine*, 5:221, 2025. DOI: 10.1038/s43856-025-00946-z. URL <https://doi.org/10.1038/s43856-025-00946-z>

a patient is likely to receive from a treatment based on their specific characteristics (**biomarkers**).

CAE: Concrete Autoencoder. RElationship between covariates and disease outcomes instead of relationship between individual treatment effects.

Modeling relationship: between covariates and disease outcomes

Prognostic model whose goal is to predict patient's likely outcome based on their baseline characteristics (*covariates*), irrespective of any treatment they may receive.

Let $X = (x_1, x_2, \dots, x_p)$ be vector of patient's baseline covariates (e.g., genetic biomarkers, age...)

Let Y be disease outcome (e.g., $Y = 1 \rightarrow$ disease progresses)

Goal: Learn f that models probability of outcome given covariates:

$$f(X) \approx P(Y = 1|X)$$

A standard logistic regression could model as:

$$\log \left(\frac{P(Y = 1|X)}{1 - P(Y = 1|X)} \right) = \beta_0 + \sum_{j=1}^p \beta_j x_j$$

CAE would be used for feature selection. In multi dimensional space where $p >>$, the goal is to identify small but relevant subset of original covariates $X_s \subset X$. Autoencoder learns to reconstruct input features using a *concrete selector* layer that is forced to choose only a few features. FINAL prognostic model is then built using ONLY selected subset:

$$f(X_s) \approx P(Y = 1|X_s)$$

Output: Model finds prognostic biomarkers that is features associated with the outcome itself.

Modeling relationship between individual treatment effects:

Predictive causal model whose goal is not just to predict outcome, but to predict how the outcome *changes* when a patient receives a specific treatment vs. control. (**Approach used by DeepRAB**).

Consider treatment variable W , where $W = 1$ for active treatment, $W = 0$ for control.

For patient with covariates X :

- Y^1 : Outcome patient will experience if they received treatment $W = 1$
- Y^0 : Same but for control $W = 0$

Individual Treatment Effect (ITE): $ITE = Y^1 - Y^0$ can only ever observe one for a given patient.

Goal is to learn model that estimates **CATE** $\tau(x)$

DeepRAB approaches such problem by learning to estimate the outcome under both scenarios. it learns 2 functions (2 heads of *NNs*)

- $\hat{\mu}_1(x) = \mathbb{E}[Y|X = x, W = 1]$ (predicted outcome if treated)
- $\hat{\mu}_0(x) = \mathbb{E}[Y|X = x, W = 0]$ (predicted outcome if controlled)
- CATE is then estimated as difference between the two predictions:

$$\hat{\tau}(x) = \hat{\mu}_1(x) - \hat{\mu}_0(x)$$

- The model is trained to minimize the prediction error on the observed outcomes (e.g., using data from a randomized clinical trial where some patients received treatment and others received control).

Core Mathematical Concepts:

Causal Inference Framework:

Neyman-Rubin potential outcome for causal inference

- Let X_i be the vector of baseline biomarkers for patient i , A_i be the treatment assignment ($A_i = 1$ for treatment, $A_i = -1$ for control), and Y_i be the observed outcome.
- The model assumes two potential outcomes for each patient: $Y_i(1)$ (outcome if treated) and $Y_i(-1)$ (outcome if on control). We only observe one of these.
- The conditional expectation of the outcome is modeled as:

$$\mathbb{E}[Y|A, X] = Z(X)A + H(X)$$

where:

- $H(X) = \frac{1}{2}[\mathbb{E}[Y|A = 1, X] + \mathbb{E}[Y|A = -1, X]]$ represents the **prognostic effect** of the biomarkers X .
- $Z(X) = \frac{1}{2}[\mathbb{E}[Y|A = 1, X] - \mathbb{E}[Y|A = -1, X]]$ is the **contrast function** that reflects the **heterogeneous treatment effect** given biomarkers X . The goal of DeepRAB is to accurately estimate this function $Z(X)$.

DeepRAB Model Architecture

Deep Neural Network (DNN) composed of 3 main components designed to estimate a *personalized benefit score*, $f(x)$, which is a monotonic transformation of the treatment effect function $Z(X)$.

- **(Component 1) Encoder / Biomarker Selection Layer:** Input layer that performs feature selection using techniques from Concrete Autoencoders (CAE). It learns to select user-specified number k of most informative biomarkers from full set of p input biomarkers. Selection is achieved by learning weight vector $\beta_j^{(0)}$ for each of the k nodes in this layer. The weights are generated using Gumbel-Softmax reparameterization (allows differentiable approximation to sampling from categorical distribution)

- Probability of selecting j -th biomarker for i -th node in this layer:

$$\beta_{ij}^{(0)} = \frac{\exp((\log \alpha_i + g_j) / T)}{\sum_{t=1}^p \exp((\log \alpha_t + g_t) / T)}$$

- α and g are learnable parameters, T is temperature parameter that is annealed towards 0 during training. As $T \rightarrow 0$, vector $\beta_j^{(0)}$ becomes a *one-hot* vector selecting a single biomarker from original input features x . Output of layer is $z^{(1)}$ which is vector of k selected biomarkers

- **(Component 2) Decoder / Hidden Layers:** Selected biomarkers $z^{(1)}$ are fed into (standard) Multi-Layer Perceptron (MLP) of $h - 1$ hidden layers that model the potentially complex and non-linear relationships between selected biomarkers and treatment effect

- Output of each hidden layer $d^{(j)}$:

$$\begin{aligned} d^{(1)} &= \phi_1(W^{(1)}z^{(1)} + b^{(1)}) \\ d^{(j)} &= \phi_j(W^{(j-1)}d^{(j-1)} + b^{(j-1)}) \end{aligned}$$

- W and b are standard weight matrices and bias vectors, ϕ is non-linear activation function

- **(Component 3) Output Layer and Loss Function:** It produces personalized benefit score $f(x)$. Model is trained by minimizing a specific loss function based on **A-learning** (Advantage-learning) which is designed to directly estimate optimal Individualized Treatment Rule (ITR) without needing to model prognostic function $H(X)$.

- A-learning loss function defined as:

$$\mathcal{F}(\theta, x_i, y_i) = \frac{1}{n} \sum_{i=1}^n M\{Y_i, (A_i - \pi(x_i))f(x_i), \theta\}$$

- θ set all trainable parameters of network
- $\pi(X) = P(A = 1|X)$ is propensity score = probability of receiving treatment given covariates X . In 1:1 randomized trial $= \pi(X) = 0.5$
- $M(u, v)$ is loss function that depends on outcome type. For continuous, it is squared error loss $M(u, v) = (u - v)^2$; for binary, it is logistic loss $M(u, v) = u \log(1 + \exp(-v))$

Model Training / Evaluation

- **Training:** θ (including α FS parameters in encoder and W, b in decoder) are optimized by min. A-learning \mathcal{F} (e.g., Adam optimizer)
- **Hyperparameter Tuning:** 10-fold CV on training via grid search
- **Evaluation Metric:** AUC. $\hat{f}(X)$ to rank patients. Vary cutoff of score to generate ROC by comparing predicted vs. true treatment (known in the simulations), then calculate AUC

Biomarker Identification:

One of the key goals of DeepRAB is to facilitate **predictive biomarker identification**. Thus after training → apply form of **model interpretability** analysis to determine which input features (biomarkers) *most strongly influence* model's prediction of high or low treatment effect $\hat{\tau}(x)$. e.g., methods:

- Gradient-based feature attribution
- Permutation feature importance
- SHAP values

Performance of mathematical framework is evaluated quantitatively using *simulated* and *real trial data*. Evaluation based on DeepRAB's ability to identify patient subgroups with enhanced treatment responses → ranking by $\hat{\tau}(x)$ separates patients who *truly* benefit from those who do not.

Mathematical Distinction

Prognostic Model (Covariates → Outcome): 'Prognostic' implies info about likely course of disease e.g., disease recurrence, progression, likelihood death. A **key attribute** of purely prognostic marker is its predictive value is *independent of the specific treatment being administered, that is, biomarker's ability to predict good/bad outcome is present in both treated and untreated*

- Y Outcome of interest (e.g., survival time, disease progression score)
- X Biomarker measurement (e.g., gene expression level, T_1 time)
- W Binary treatment indicator

$$\mathbb{E}[Y|X, W] = \beta_0 + \beta_X X + \beta_W W + \beta_{XW}(X \cdot W)$$

Conditions:

1. Must be associated with outcome: Coefficient for biomarker itself must be significant

$$\beta_X \neq 0$$

X provides information about Y even when $W = 0$

2. Effect must NOT depend on treatment: No significant interaction between biomarker and treatment

$$\beta_{XW} = 0$$

Effect of X on Y is same for both $W = 1$ and $W = 0$. Graphically, the lines representing the relationship between X and Y for the treated and control groups are **parallel**. Here, treatment effect β_W is constant $\forall x_i \in X$

- Models $P(Y|X)$
- '*Given x patient's biomarker, what is their likely prognosis?*'
- Biomarkers found: Prognostic. Such features predict outcome regardless of treatment

Predictive Model of Treatment Effect (Covariates → Treatment Effect):

- Models $\mathbb{E}[Y^{(1)} - Y^{(0)}|X = x]$
- '*Given patient's biomarkers, how much benefit will they get from treatment vs. control?*'
- Biomarkers found: Predictive. Outcome prediction & Prediction of response difference to treatment. e.g., X biomarker might not have relationship with outcome in control but a strong relationship in treated

Summary

Subgroup identification and modeling treatment effect heterogeneity with predictive biomarker identification (feature selection method) being key component and outcome of process.

June 12th / 2025

Variant Effect Predictors (VEPS)

Computational tools whose primary function is to assess potential functional impact of genetic variants (particularly **missense**) which *cause a change in AA sequence of a protein.* → crucial in addressing variants of unknown significance VUS.

Foundational Principles

Evolutionary Conservation, Sequence Homology, and Structural Information

- **SIFT (Sorting Intolerant From Tolerant):** Operates on principle that important AA will be conserved throughout evolution. It predicts whether AA substitution will impact protein function by analyzing conservation across multiple species.
- **PolyPhen-2 (Polymorphism Phenotyping v2):** Predictor takes multifaceted approach. Evaluates physicochemical differences between AA, position of substitution within protein's structure, proximity to functional domains on top of evolutionary conservation.
- **CADD (combined annotation dependent depletion):** Scores deleteriousness of variants by integrating multiple annotations. Trained by comparing a vast *set* of observed human variations (mostly neutral) against simulated mutations to learn how to distinguish between them. Such approach allows to score coding and non-coding variants

MODERN AI-BASED VEPs use *transformers – based* architecture which allows model to weigh importance of different parts of the sequence to **understand context.** Thus it aims to learn fundamental language and rules of protein structure and function *without being told* which variants are pathogenic.

ACCURATE PROTEOME-WIDE MISSENSE VARIANT EFFECT PREDICTION WITH ALPHAMISSENSE^{3 4}

α-Missense

By Google DeepMind. It predicts pathogenicity of missense variants. Architectures is inspired by *Evoformer* block used in α -fold

Core idea is to iteratively refine 2 key representations:

- **MSA representation:** Captures evolutionary information
- **Pair representation:** Captures spatial and relational information between AA pairs

Input Representation:

- **Target Sequence:** Primary protein sequence of length L is typically one-hot encoded into matrix $S_{\text{target}} \in \{0, 1\}^{L \times 20}$, where 20 is AAs
- **MSA Representation (M):** N aligned sequences of length L represented a tensor $M \in \mathbb{R}^{N \times L \times d_{msa}}$. Each position (i, j) in tensor is an embedding for AA at residue j in sequence i of alignment
- **Pair Representation (P):** Tensor $P \in \mathbb{R}^{L \times L \times d_{pair}}$ is initialized to store information about pairs of residues (i, j) . Can be initialized with info about relative positions of residues in the sequence, i.e., $j - i$.

Pair Representation: Tensor built and maintained by *Evoformer*. It stores and refine model's understanding of the relationship between every pair of residues in the protein.

Tensor P of dimensions $L * L * d_{pair}$ where L length protein sequence and d_{pair} is # features the model stores for each pair. *Just like $L \times L$ matrix but instead of scalar in ij , a high dimensional vector of features describing relationship between ij*

Evoformer Block

The model consists of series of stacked *Evoformer* blocks. Each block takes M and P and outputs updated M' and P' . Info is allowed to flow back and forth between MSA (evolutionary context) and implicit structural context (pair representation)

The Block:

³ Jun Cheng, Guido Novati, Joshua Pan, Clare Bycroft, Akvilė Žemgalytė, Taylor Applebaum, Alexander Pritzel, Lai Hong Wong, Michal Zielinski, Tobias Sargeant, Rosalia G. Schneider, Andrew W. Senior, John Jumper, Demis Hassabis, Pushmeet Kohli, and Žiga Avsec. Accurate proteome-wide missense variant effect prediction with alphamissense. *Science*, 381(6664):eadg7492, 2023.
DOI: [10.1126/science.adg7492](https://doi.org/10.1126/science.adg7492). URL <https://www.science.org/doi/abs/10.1126/science.adg7492>

⁴ J. Jumper, R. Evans, A. Pritzel, and et al. Highly accurate protein structure prediction with alphafold. *Nature*, 596(7873), 2021

MSA Representation Update (with Axial Attention) MSA representation M is updated using axial attention which is applied independently along the rows and columns.

Row-wise Attention (within each sequence): For each sequence i in MSA i.e. each row of M , standard self-attention mechanism is applied across residues $j = 1, \dots, L \rightarrow$ Model learns relationship between different residues within the same sequence. For a single row:

$$\text{Attention}(Q, K, V) = \text{softmax} \left(\frac{QK^T}{\sqrt{d_k}} \right) V \quad (1)$$

where Q, K, V = linear projections of row's residue embeddings, augmented with information from the pair representation P . MSA representation is updated to M_{row}

Column-wise Attention (across sequences): A second attention mechanism is applied to each column j of intermediate MSA M_{row} Model then learns which sequences in the alignment are most informative for x residue position. Capable of identifying highly conserved positions and co-evolving mutations across different sequences (signal for functional importance).

Communication: MSA to Pair Representation

Information from updated MSA representation is used to update pair representation. Crucial step where *evolutionary information informs spatial relationships*.

Some context:

Note that pair representation has been used in 2 places within the single block up to this point. This is the core cyclical logic of the Evoformer architecture → **Iterative refinement**(reasoning cycle). Model uses *current* pair representation P to *bias/guide* attention mechanism. Thus instead of e.g., $Scores = QK^T$, it becomes $Scores = QK^T + Bias_{ij}$. A large $Bias_{ij}$ = model is forced to focus on relationship between residues i and j when updating MSA representation.

- Update is often achieved using \otimes -like operation on columns of MSA representation (correlation matrix). For a pair of residues (i, j) model takes corresponding columns from MSA embedding, $M_{:,i}$ and $M_{:,j}$, and combines them. Simplified view of update for pair (i, j) :

$$P'_{ij} = \text{LinearLayer} \left(\sum_{k=1}^N (W_1 M_{ki}) \otimes (W_2 M_{kj}) \right) \quad (2)$$

$M_{:,i}$ is AA at **residue position** i across all the different sequences in the alignment. Embedding captures evolutionary variation at that specific site in the protein, thus operation computes covariance matrix over MSA embeddings for each pair of residues **columns** $M_{:,i}$ and $M_{:,j}$ capturing co-evolutionary signals

Pair Representation Update: Tensor P undergoes own refinement using series of convolutional layers or axial attention layers or axial attention layers applied over $L * L$ grid. Analogous to refining distance so model is allowed to enforce geometric consistency rules like e.g., triangle inequality, on the relationship between residues. P' is fed back into MSA update in the next Evoformer block.

Prediction Head and Pathogenicity Score

To make a prediction for a x missense variant (e.g., wild type AA a_{wt} at i is replaced by a_{mut}) a prediction head is used.

- Extract final embedding for residue i from target sequence's representation $h_i \in \mathbb{R}^{d_{model}}$
- During self-supervised pre-training, a classification head (e.g., linear layer followed by softmax) is used to predict probability distribution over all 20 AAs for x position

$$\text{Logits} = W_{pretrain} h_i + b_{pretrain} \quad (3)$$

$$P(\text{amino acid}_j | \text{context}) = \text{softmax}(\text{Logits})_j = \frac{e^{\logit_j}}{\sum_{k=1}^{20} e^{\logit_k}} \quad (4)$$

For final pathogenicity prediction, fine-tuning process trains simpler head. Head takes final representation h_i (which contains information about a_{wt} and a_{mut} context) and projects it to a single scalar value. This is a **binary classification task**.

$$\alpha = \sigma(W_{patho} h_i + b_{patho}) \quad (5)$$

α is final pathogenicity score

SIFT: PREDICTING AMINO ACID CHANGES THAT AFFECT PROTEIN FUNCTION⁵

⁵ P. C. Ng and S. Henikoff. Sift: predicting amino acid changes that affect protein function. *Nucleic Acids Research*, 31(13):3812–3814, 2003. DOI: 10.1093/nar/gkg509

SIFT Algorithm

Evolutionary conservation. It says that AA positions critical for a protein's structure/function will be conserved across homologous sequences from different species. Therefore, a substitution at a highly conserved position is likely to be deleterious, whereas a substitution at a highly variable position is more likely to be tolerated without significant functional consequence. SIFT quantifies this by calculating a score based on the probability of observing a particular A at a specific position, derived from a multiple sequence alignment (MSA).

Step 1: Multiple Sequence Alignment (MSA) Generation

Given a query protein sequence ,the first step is to gather set of homologous sequences from large protein database (e.g., Swiss-Prot/TrEMBL) using an algorithm like *PSI – BLAST* (Position-Specific Iterated Basic Local Alignment Search Tool). Result is an MSA where related sequences are aligned, revealing *patterns of conservation and variation at each position*.

Step 2: Position-Specific Probability Matrix (PSPM) Calculation

Core calculation in SIFT. For each position i in the alignment, the algorithm computes a **probability distribution** over the 20 AAs

Let P_{ij} be probability of AA j occurring at position i . Probabilities are calculated from the **observed frequencies** of AAs at that position in the MSA. To handle sampling bias (e.g., over-representation of very similar sequences) and zero-frequency events (amino acids not seen at a position), a Bayesian approach using a Dirichlet mixture as prior probabilities is used

Simplified representation using pseudocounts:

$$P_{ij} = \frac{n_{ij} + b_j}{N_i + B} \quad (6)$$

- n_{ij} is weighted count of sequences in the MSA having AA j at position i . Sequence weights are used to down-weight redundant, highly similar sequences
- $N_i = \sum_{j=1}^{20} n_{ij}$ is total weighted count of sequences at position i
- b_j is the pseudocount for AA j . Derived from a prior probability distribution (substitution matrix e.g., BLOSUMXX)
- $B = \sum_{j=1}^{20} b_j$ is the total number of pseudocounts

The set of these probabilities for a given position i , $\{P_{i1}, P_{i2}, \dots, P_{i20}\}$, forms the Position-Specific Probability Matrix (a vector for position i)

and satisfies:

$$\sum_{j=1}^{20} P_{ij} = 1 \quad (7)$$

Step 3: SIFT Score Calculation and Classification

The SIFT score for a given substitution from the wild-type AA (aa_{wt}) to a mutant AA (aa_{mut}) at position i is the probability of observing that mutant AA at that position, as derived from the PSPM.

$$\text{SIFT Score}(i, aa_{wt} \rightarrow aa_{mut}) = P_{i,aa_{mut}} \quad (8)$$

This score is a *measure of tolerance*. A high score (high probability) indicates that the substitution is commonly observed in homologs and is therefore predicted to be tolerated. A low score indicates the substitution is rare and likely not tolerated (deleterious).

The final classification is made by applying a threshold, typically 0.05:

$$\text{Prediction} = \begin{cases} \text{Deleterious} & \text{if SIFT Score} < 0.05 \\ \text{Tolerated} & \text{if SIFT Score} \geq 0.05 \end{cases} \quad (9)$$

Step 4: Conservation Index

SIFT's predictions are more reliable for positions that are *highly conserved*. To **quantify**, a conservation index is calculated for each position i , which is derived from the information content or negative of Shannon's Entropy of the probability distribution P_i .

First, the Shannon's Entropy H_i for position i is calculated:

$$H_i = - \sum_{j=1}^{20} P_{ij} \log_2(P_{ij}) \quad (10)$$

The entropy H_i is a measure of uncertainty/variability at position i . It ranges from 0 (perfect conservation, only 1 AA is possible) to $\log_2(20)$ (all 20 AA equally likely).

The Conservation Index C_i is then the information content: difference between the maximum possible entropy and the observed entropy:

$$C_i = \log_2(20) - H_i \quad (11)$$

A high conservation index ($C_i \rightarrow \log_2(20)$) indicates low entropy and high conservation, making SIFT prediction at that site more reliable. Likewise $C_i \rightarrow 0$ indicates high variability and predictions at such sites are considered less certain.

DNABERT: PRE-TRAINED BIDIRECTIONAL ENCODER REPRESENTATIONS FROM TRANSFORMERS MODEL FOR DNA-LANGUAGE IN GENOME⁶

Methods to improve statistical power of gene-level association tests by partitioning rare variants into K functional categories (S_1, \dots, S_K). Methods outperform standard tests that treat all variants equally especially when different functional categories have different magnitudes.

Method 1: Omnibus SKAT (oSKAT)

Performs separate Sequence Kernel Association Test (SKAT) for each of K variant sets yielding K p-values (p_1, \dots, p_K). p-values are combined using *Simes' method* to produce single gene-level p-value (then adjust for correlation between tests)

Simes' p-value:

$$p_{\text{Simes}} = \min_{i=1, \dots, K} \frac{K \cdot p_{(i)}}{i} \quad (12)$$

where $p_{(i)}$ is i -th smallest p-value.

Method 2: Functional SKAT (F-SKAT)

F-SKAT = unified variance component test within single mixed model. Overall genetic effect for an individual, γ_i :

$$\gamma_i = \sum_{k=1}^K \gamma_{ik}, \quad \text{where } \gamma_{ik} = \sum_{j \in S_k} w_j \beta_j G_{ij} \quad (13)$$

Assumes variant effects β_j for each category k are random variables drawn from distribution with a category-specific variance component, τ_k :

$$\beta_j \sim N(0, \tau_k) \quad \text{for } j \in S_k \quad (14)$$

H_0 is 'all variance components are zero' = no genetic effect from any category:

$$H_0 : \tau_1 = \tau_2 = \dots = \tau_K = 0 \quad (15)$$

The F-SKAT score test statistic is an optimal linear combination of the individual SKAT statistics (Q_k) for each category:

$$Q_F = \sum_{k=1}^K \lambda_k Q_k \quad (16)$$

This statistic follows mixture of χ^2 distributions from which p-value can be derived.

Simulations and real data analyses show F-SKAT is generally the most powerful and flexible approach.

⁶ Yanrong Ji, Zhihan Zhou, Han Liu, and Ramana V Davuluri. Dnabert: pre-trained bidirectional encoder representations from transformers model for dna-language in genome. *Bioinformatics*, 37(15):2112–2120, 02 2021. ISSN 1367-4803. doi: 10.1093/bioinformatics/btab083. URL <https://doi.org/10.1093/bioinformatics/btab083>

June 13th /2025

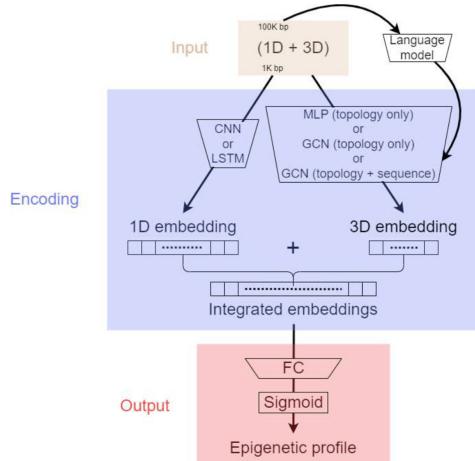
MULTIMODAL LEARNING OF NONCODING VARIANT EFFECTS USING GENOME SEQUENCE AND CHROMATIN STRUCTURE ⁷

- DL framework to predict the effects of noncoding genetic variants by integrating 1D local genome sequence and 3D global chromatin structure. Mathematical architecture overcomes challenges of multimodal data integration (e.g., significance difference in data resolution)

1D Sequence and Epigenetic data from DeepSEA (already processed). $5.2 * 10^6$ samples where each sample is 1kb DNA sequence. Each sequence is labeled with 919 binary values corresponding to epigenetic effects across 148 cell lines

3D Structure data was sourced from Hi-C experiments on the ENCODE portal which provides genome-wide chromatin interaction frequency matrices for cell lines GM12878, IMR90, and K562. Matrices represent 3D proximity of different genomic regions (100kb resolution).

Mathematical Framework



⁷ Wuwei Tan and Yang Shen. Multimodal learning of noncoding variant effects using genome sequence and chromatin structure. *Bioinformatics*, 39(9):btad541, 09 2023. ISSN 1367-4811. DOI: 10.1093/bioinformatics/btad541. URL <https://doi.org/10.1093/bioinformatics/btad541>

$1D$ Sequence Encoding: CNNs and RNNs ((biLSTM))

For a $1D$ input sequence S and kernel K of size k , the output feature map C at position i :

$$C_i = f\left(\sum_{j=1}^k K_j * S_{i+j-1} + b\right)$$

f is non-linear

and,

$$h_{t-1} : h_t = \tanh(W_{hh}H_{t-1} + W_{xh}x_t + b_h)$$

$3D$ Structure Encoding: GNNs

- GCNs capture complex, non-grid-like topology of chromatin interactions
- Feature vector (embedding) h_i^{l+1} for node i at layer $l + 1$

$$H^{(l+1)} = \sigma(\hat{D}^{-\frac{1}{2}} \hat{A} \hat{D}^{-\frac{1}{2}} H^{(l)} W^{(l)})$$

where $H^{(l)}$ is matrix of node features at layer l , \hat{D} is diagonal degree matrix of \hat{A} and $\hat{A} = A + I$ (adjacency matrix).

Pre-trained DNA language model used: m **DNABERT**.

Training | Prediction

Loss: Binary cross-entropy. For single data sample:

$$L = - \sum_{i=1}^{919} [y_i \log(\hat{y}_i) + (1 - y_i) \log(1 - \hat{y}_i)]$$

Regularization: L1 and L2

L1 penalty: $\lambda_1 \sum |w|$

L2 penalty: $\lambda_2 \sum w^2$

Results

Sequence-only models: DeepSEA, DanQ, Sei-based model.

Incorporating $3D$ chromatin structure into model outperforms sequence-only. GCN + DNABERT for structure embedding → Improved AUPRC by over 10%.

Performance Variant Effect Prediction:

- **eQTL pred.:** Structure-informed models as feature generators yielded better prediction of noncoding variant effects on gene expression vs. sequence-only

- **Pathogenicity pred.:** in *Unsupervised* setting, AUROC ≈ 0.75 and AUPRC ≈ 0.15. In *Supervised* setting (even with 100 labeled training samples) AUROC > 0.8 and AUPRC > 0.25

Findings

- 3D chromatin structure helps explain disparity between DNA sequence similarity and epigenetic profile similarity. DNA regions far apart in the 1D sequence *but have* similar epigenetic profiles tend to be closer in 3D space
- Models corrected bias of sequence-only predictors (where sequence similarity was poor indicator of epigenetic similarity)
- Using GCNs to embed 3D structure data was **key strategy**. Chromatin interaction modeled as graphs.
- Structure-informed models identified 2 motifs related to long-range interactions - POU3F1 and TFDP1) which were missed by DanQ BUT at the same time they missed 7 motifs that DanQ found because they lacked **long-range interaction patters**
- Little to no feature engineering. applicable to multiple types of mutations (insertions, deletions ...)
- DID NOT outperformed CADD in their OWN but when combined
 - 7 motifs missed that were found by sequence-only DanQ likely because motifs lacked long-range interactions that Hi-C data captures. *Using higher-resolution chromatin structure data e.g., from Micro-C experiments, may help.*

STANDARDS AND GUIDELINES FOR THE INTERPRETATION OF SEQUENCE VARIANTS: A JOINT CONSENSUS RECOMMENDATION OF THE AMERICAN COLLEGE OF MEDICAL GENETICS AND GENOMICS AND THE ASSOCIATION FOR MOLECULAR PATHOLOGY⁸

- Framework for standardizing interpretation and classification of genetic variants for Mendelian diseases.

Goal: To create robust, evidence-based system for consistent variant classification in a clinical context.

5-Tier Classification System

- **Pathogenic**

⁸ S Richards, N Aziz, S Bale, D Bick, S Das, J Gastier-Foster, W W Grody, M Hegde, E Lyon, E Spector, K Voelkerding, H L Rehm, and ACMG Laboratory Quality Assurance Committee. Standards and guidelines for the interpretation of sequence variants: a joint consensus recommendation of the american college of medical genetics and genomics and the association for molecular pathology. *Genetics in Medicine*, 17(5):405–424, 2015. DOI: 10.1038/gim.2015.30

- **Likely Pathogenic****
- **Uncertain Significance (VUS)**
- **Likely Benign****
- **Benign**

** Represent > 90%

A Variant of Uncertain Significance (VUS) should not be used in clinical decision-making and efforts should be made to resolve its classification

Evidence Framework

For combining different types of evidence to arrive at x classification. Authors chose a rule-based system for combining evidence codes instead of numerical scoring system. Exact numerical relationships cannot be proved yet. Given that the impact of a piece of evidence is often context-dependent they would fail to capture context.

Pathogenic Evidence Strengths:

- **PVS1 (very strong)** Predicted Null variant (e.g., nonsense, frameshift) in a gene where loss of function is a known mechanism of
- **PS1-PS4 (strong)** Includes evidence like a variant being previously established pathogenic missense variant, confirmed de novo variant in patient with consistent phenotype or variant showing statistically significant increased prevalence in affected individuals vs. controls
- **PM1-PM6 (Moderate)** Includes evidence e.g., being located in a mutational hotspot, being absent from controls in population databases (e.g., ExAC, 1000 Genomes) or being detected *in trans* with another pathogenic variant for a recessive disorder
- **PP1-PP5 (Supporting)** Evidence like co-segregation with disease in a family, multiple lines of computational evidence (tools like SIFT, PolyPhen-2) supporting a damaging effect or a patient's phenotype being highly specific for the gene

Benign Evidence Strengths

- **BA1(Stand-Alone):** Allele frequency is > 5% in a large population database

- **BS1-BS4 (Strong):** Includes evidence like having an allele frequency greater than expected for x disorder, being observed in a healthy adult for a fully penetrant childhood disease, well established functional studies showing no damaging effect.
- **Supporting:** Includes evidence like being a missense variant in a gene where only truncating variants are known to cause disease or multiple computational tools suggesting no impact

Classification Rules | 'Algorithm Classifier'

Pathogenic

1. 1 Very Strong (PVS₁) **AND** one of the following:
 - (a) ≥ 1 Strong (PS₁–PS₄)
 - (b) ≥ 2 Moderate (PM₁–PM₆)
 - (c) 1 Moderate (PM₁–PM₆) **AND** 1 Supporting (PP₁–PP₅)
 - (d) ≥ 2 Supporting (PP₁–PP₅)
2. ≥ 2 Strong (PS₁–PS₄)
3. 1 Strong (PS₁–PS₄) **AND** one of the following:
 - (a) ≥ 3 Moderate (PM₁–PM₆)
 - (b) 2 Moderate (PM₁–PM₆) **AND** ≥ 2 Supporting (PP₁–PP₅)
 - (c) 1 Moderate (PM₁–PM₆) **AND** ≥ 4 Supporting (PP₁–PP₅)

Table 1: Rules for Combining Criteria to Classify Sequence Variants

Likely Pathogenic

1. 1 Very Strong (PVS₁) **AND** 1 Moderate (PM₁–PM₆)
2. 1 Strong (PS₁–PS₄) **AND** 1–2 Moderate (PM₁–PM₆)
3. 1 Strong (PS₁–PS₄) **AND** ≥ 2 Supporting (PP₁–PP₅)
4. ≥ 3 Moderate (PM₁–PM₆)
5. 2 Moderate (PM₁–PM₆) **AND** ≥ 2 Supporting (PP₁–PP₅)
6. 1 Moderate (PM₁–PM₆) **AND** ≥ 4 Supporting (PP₁–PP₅)

Benign

1. 1 Stand-Alone (BA₁) **OR**
2. ≥ 2 Strong (BS₁–BS₄)

Likely Benign

1. 1 Strong (BS₁–BS₄) **AND** 1 Supporting (BP₁–BP₇) **OR**
2. ≥ 2 Supporting (BP₁–BP₇)

Variants should be classified as Uncertain Significance if other criteria are unmet or the criteria for benign and pathogenic are contradictory.

June 15th /2025

► ACMG/AMP provides standardized system for variant classification but it is **not** a rigid algorithm. The guidelines are built with flexibility and state that expert judgement is required.

► Exceptions to guidelines may be cases where context, gene-specific knowledge, whether quality of evidence allows a rule to be applied differently or not.

Guidelines instead of rigid point-based system

Authors state:

'that the assignment of specific points for each criterion implied a level of quantitative understanding...that is currently not supported scientifically and does not take into account the complexity of interpreting genetic evidence'

Entire framework is thus built on the premise of *allowing expert curation and judgement* to upgrade/downgrade/ignore a piece of evidence based on context.

Exceptions and Context-Dependent Rule ⁹

- **Context-Dependent Evidence Strength:** The weight of several criteria can be adjusted based on the available data. For example:
 - The evidence for a variant co-segregating with disease in a family (**PP1**) can be upgraded from 'Supporting' to 'Moderate' or 'Strong' if data from multiple large families is available
 - Observing a variant *in trans* with a known pathogenic variant for a recessive disorder (**PM3**) can be upgraded from 'Moderate' to 'Strong' if this is observed multiple times
- **Gene-Specific Disease Mechanisms:** The applicability of certain rules depends entirely on the known biology of the gene in question.

⁹ S Richards, N Aziz, S Bale, D Bick, S Das, J Gastier-Foster, W W Grody, M Hegde, E Lyon, E Spector, K Voelkerding, H L Rehm, and ACMG Laboratory Quality Assurance Committee. Standards and guidelines for the interpretation of sequence variants: a joint consensus recommendation of the american college of medical genetics and genomics and the association for molecular pathology. *Genetics in Medicine*, 17(5):405–424, 2015. DOI: 10.1038/gim.2015.30

- The **PVS1** (predicted null variant) criterion is considered 'Very Strong' pathogenic evidence, but only for genes where loss-of-function is a known disease mechanism. For many cardiovascular genes like *MYH7*, where missense variants are the primary cause of disease, a heterozygous null variant is not pathogenic, and PVS1 does not apply
- **Evolving Guidelines:** The framework is intended to evolve as new data becomes available.
 - The **PM2** criterion (variant is absent from controls in population databases) was originally weighted as 'Moderate'. However, subsequent analysis by the ClinGen consortium recommended downgrading its strength to 'Supporting' on its own, acknowledging that most very rare or novel variants are benign

June 16th /2025

CARDIOVASCULAR DISEASES AND THE ROLE OF MISSENSE VARIANTS

Missense Variant

single nucleotide change in DNA resulting in another AA incorporated in *x* protein. Possible outcome → *benign* or *severely disruptive* which often leads to alteration of protein structure and thus function.

In cardiogenetics, missense variants are a common mechanism of disease particularly for *inherited heart muscle and rhythm disorder*.

Key examples:

Hypertrophic Cardiomyopathy (HCM): Thickening of the heart muscle. **MYH7** and **MYBPC3** are the most common involved genes which encode proteins of the sarcomere (heart's contractile unit). The variants often lead to an altered protein that gets incorporated into the sarcomere and disrupts its function often through a **Dominant-Negative Effect** where the abnormal protein interferes with the function of the normal protein from the other allele.

► As noted in ACMG/AMP, simple loss-of-function (null) variants in many of such genes are more likely to be benign. Therefore, presence of *x* faulty component (missense variant) in the machine contributes to its abnormal behaviour *not* the lack of such component.

Dilated Cardiomyopathy (DCM): Enlarged and weakened left ventricle. Missense variants in genes encoding proteins of the sarcomere (e.g., **TTN**), cytoskeleton, nuclear envelop are major the major cause *but* loss-of-function variants can also cause DCM. These variants can compromise the structural integrity and force-generating capacity of the heart muscle cells

Arrhythmogenic Cardiomyopathy (ACM): Tissue is replaced by fatty and fibrous tissue which lead to arrhythmias. Often caused by

missense variants in genes that encode **desmosomal** protein (e.g., PKP2, DSG2, DSP) which are essential for holding heart cells together. A single AA change can disrupt cell-cell junctions → cell death / disease.

Cardiac Channelopathies: Group of disorders caused by mutations in genes encoding cardiac ion channels that control the heart's electrical activity. Their primary mechanism is often missense variants.

GENETICS OF MYOCARDIAL INTERSTITIAL FIBROSIS IN THE HUMAN HEART AND ASSOCIATION WITH DISEASE¹⁰

Myocardial Interstitial Fibrosis:

Pathological scarring of heart muscle. Excessive deposition of extracellular matrix proteins (e.g., collagen) that lead to cardiac stiffness, impaired function, and ultimately, heart failure.

Monogenetic Causes:

These are typically rare variants with a large effect size

Sarcomeric Genes in Hypertrophic Cardiomyopathy (HCM) Most common genetic heart disease. A major cause are mutations that code for the sarcomere (contractile unit).

Frequently implicated genes:

- **MYH7 (Myosin Heavy Chain 7):** Encodes β -myosin heavy chain (motor protein for muscle contraction)
- **MYBPC3 (Myosin Binding Protein C, Cardiac):** Encodes protein that regulate the contraction and relaxation of the heart muscle

Other sarcomere-related genes are TNNT2, TNNI3. These mutations are thought to initiate a cascade that results in fibrosis. The presence of such mutations is strongly associated with a greater extent of myocardial fibrosis

Polygenic Risk: Cumulative Effect

For a large portion of the population the risk of developing myocardial fibrosis is not tied to a single faulty gene but rather to the combined influence of many common genetic variants. GWAS have

¹⁰ V Nauffal, P Di Achille, M D R Klarqvist, and et al. Genetics of myocardial interstitial fibrosis in the human heart and association with disease. *Nature Genetics*, 55:777–786, 2023. DOI: 10.1038/s41588-023-01371-5

identified loci associated with an increased risk of fibrosis. Such variants are linked to a variety of biological processes that when altered can promote a pro-fibrotic environment in the heart.

- **Glucose Transpose:** Variants in **SLC2A12** may alter energy metabolism within the heart
- **Iron Homeostasis:** **HFE** and **TMPRSS6** suggest a link between iron regulation and cardiac scarring
- **Oxidative Stress:** **CAMK2D** variant is involved in calcium signaling which is critical for cardiomyocyte function and survival
- **Tissue Repair and Remodelings:** **ADAMTSL1** and **VEGFC** play roles in how the extracellular matrix is maintained and repaired
- **Chromatin Remodeling:** **SMARCB1** variants have been linked to an increased risk of fibrosis. Decreased expression of it can lead to an exaggerated response to fibrotic stimuli.
 - Regardless of the trigger (monogenic or polygenic), multiple signaling pathways are implicated in downstream process of fibrosis. Genetic variations tend to converge on such pathways → Amplification Fibrotic Response.

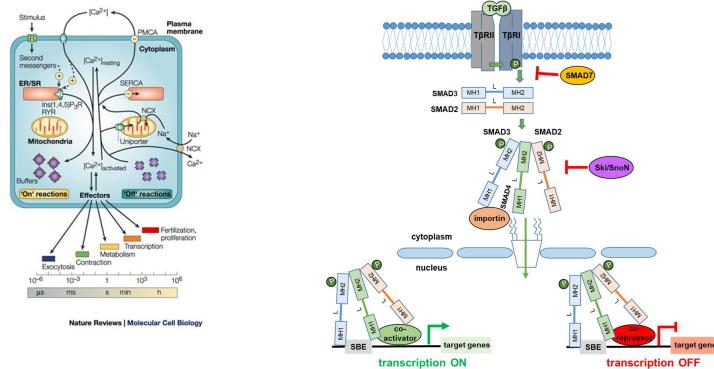


Figure 1: Signaling pathways: Calcium and TGF

Authors goal is to overcome limitation in understanding genetic basis of **myocardial interstitial fibrosis** using ML to quantify fibrosis from *cardiac magnetic resonance imaging (CMRI)* (UK biobank cohort) then identify novel genetic pathways in the disease.

U-Net with DenseNet-121 encoder

Symmetric encoder-decoder structure. Designed to capture context and localization

DenseNet-121 Encoder

Input image → hierarchical feature representations at different spatial scales.

- **Composite Layer Function (H_l)**

Let x_{l-1} output preceding layer: $x'_l = \text{Conv}(\text{ReLU}(\text{BN}(x_{l-1})))$

- **Batch Normalization (BN)** normalizes input z , scales, shifts it.
For input feature map z :

$$\text{BN}(z) = \gamma \left(\frac{z - \mu_{batch}}{\sqrt{\sigma_{batch}^2 + \epsilon}} \right) + \beta$$

where μ and σ^2 pertain to mini-batch. β is shifting parameter, γ is learning rate

- **Dense Connectivity**

Output of l is $x_l = H_l([x_0, x_1, \dots, x_{l-1}])$ composite function H_l is applied to **concatenation** of feature maps along channel dimension (e.g., $k_0 + (l+1) * k$ channels for l)

- **Transition Layers**

[Between dense blocks] Used to control complexity and *downsample spatial dimensions of feature maps*.

- 1×1 convolution to reduce number of feature maps
- 2×2 avg. pooling layer to halve height/width of feature map

$$x_{trans} = \text{AvgPool}(\text{Conv}_{1 \times 1}(x))$$

U-Net Decoder and Skip Connections

Overall framework for segmentation. Input is output from encoder.

- **Encoder-Decoder Symmetry** The U-Net has 2 paths:

- **Encoder:** DenseNet-121 **output** from different levels: e_1, e_2, e_3, e_4
- The **Decoder** symmetrically reconstructs the spatial resolution to produce segmentation map

- **Upsampling (Transposed Convolution (or deconvolution))**

let d_i be the feature map at level i in the decoder. The unsampled feature map u_i : $u_i = \text{TransConv}(d_i)$

$$L(p, g) = \alpha \cdot L_{\text{focal}} + (1 - \alpha) \cdot L_{\text{dice}} \quad (17)$$

$$L_{\text{focal}} = -(1 - p_t)^\gamma \log(p_t) \quad (18)$$

$$p_t = \begin{cases} p & \text{if } g = 1 \\ 1 - p & \text{otherwise} \end{cases}$$

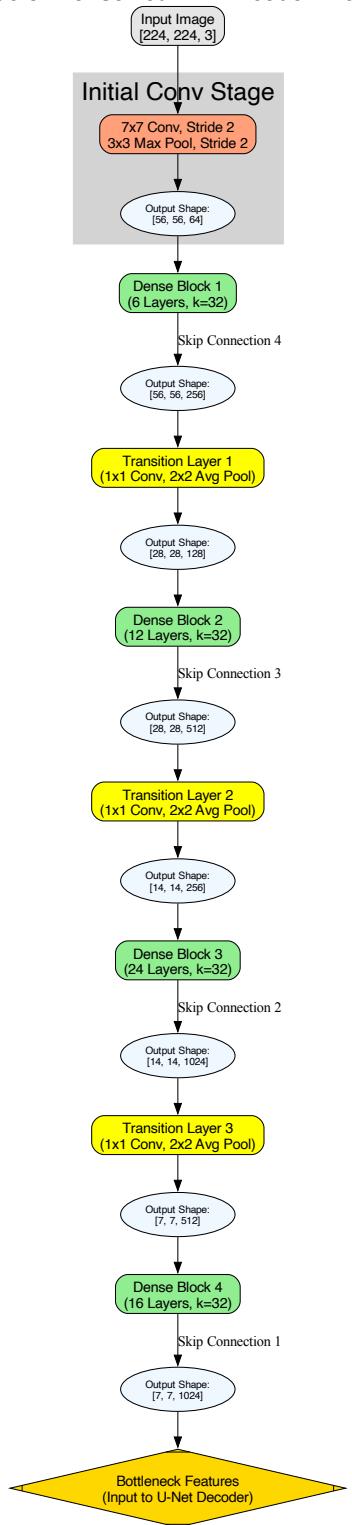
$$L_{\text{dice}} = 1 - \frac{2 \sum_i (p_i \cdot g_i) + \epsilon}{\sum_i p_i^2 + \sum_i g_i^2 + \epsilon} \quad (19)$$

PLA¹¹

¹¹

Figure 2: Model structure.

Schematic of DenseNet-121 Encoder Architecture



June 19th / 2025

Next-Generation Sequencing

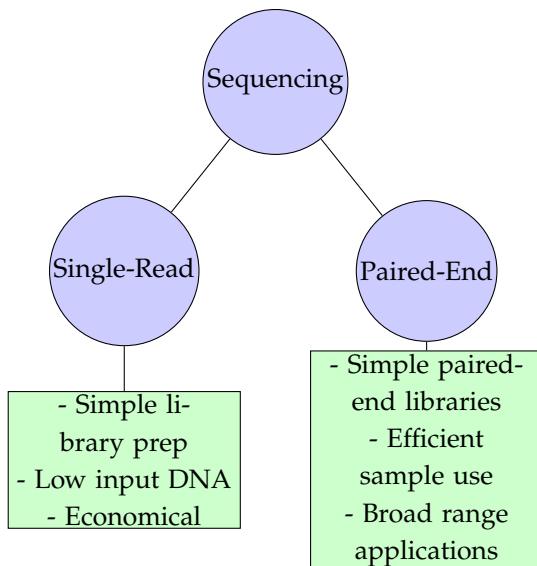
¹²

- **Read:** Single sequence produced from a sequencer.
- **Library:** Collection of DNA fragments that have been prepared for sequencing.
- **Flowcell:** Chip on which DNA is loaded and provided to the sequencer.
- **Lane:** Open portion of a flowcell. Usually for technical replicates or different samples.
- **Run:** Entire sequencing reaction from start to finish.

¹² NYU Langone Health. Next-generation sequencing analysis resources, 2023. URL <https://learn.gencore.bio.nyu.edu>. Accessed on 2025-05-30

Steps:

1. Sample collection/preparation
2. Amplification
3. Basecalling



File Formats in Genomics

FastA Format:

Most basic for reporting a sequence. Contains sequence name, description.

```
>Chr1 CHROMOSOME dumped from ADB:  
Jun/20/09 14:53; last updated: 2009-02-02  
CCCTAACCTAAACCTAAACCTAAACCTC  
TGAATCCTTAATCCCTAAATCCCTAAATCTT  
AAATCCTACATCCAT
```

DB query tools such as **blast** and **multiple-sequence alignment** programs accept only FastA. Reference Genomes are often delivered in this format.

FastQ Format:

Most widely used in sequence analysis. Output delivered from a sequencer. More information is contained in this format.

```
@SEQ_ID  
GATTGGGGTTCAAAGCAGTATCGATCAAATAGTAAATCCAT  
TTGTTCAACTCACAGTT  
+  
!''*( ((****) )%%++ )%%% . 1***-+*'')**5  
5CCF>>>>CCCCCCC65
```

Quality value characters. Lowest: ! and highest: ~:

```
! "#$%&' ()*+, -./0123456789: ;<=>?@ABCDEFGHIJKLMN
OPQRSTUWXYZ[\]^_`abcdefghijklmnopqrstuvwxyz
wxyz{|}~
```

- Sequence Header:
 - @ is sequence identifier
 - rest is sequence description
- Second line is sequence
- Third line starts with + and can have same sequence identifier appended
- Fourth are quality scores (ASCII encoded)

Sequence Header contains: Instrument name, run ID, flowcell ID, flowcell lane, tile number, X and Y coordinates of cluster, member of a pair, filtered status, control bits, index sequence, etc.

Nearly everything works with FastQ except for: Blast, Multiple sequence alignment (typically require FastA), any reference sequence (usually FastA).

Quality Scores/Q-score:

Integer value assigned to each nucleotide base. Represents estimated probability (P) that base call is incorrect. Logarithmic scale: $Q = -10 \log_{10}(P)$. **Higher Q-score** = Higher confidence = lower P (error) = higher accuracy. **Lower Q-score** = Lower confidence = higher P (error).

Examples:

- **Q10:** $P = 0.1$ (1 in 10 error). Accuracy = 90%.
- **Q20:** $P = 0.01$ (1 in 100 error). Accuracy = 99%. Often a minimum acceptable quality.
- **Q30:** $P = 0.001$ (1 in 1,000 error). Accuracy = 99.9%. Benchmark for high-quality.
- **Q40:** $P = 0.0001$ (1 in 10,000 error). Accuracy = 99.99%.

Common uses: filter bases/reads if a threshold is not met. Main purpose: provide evidence that the sequence, alignment, assembly, SNP are real and not sequencing artifacts.

SAM Format:

Sequence Alignment/Map. Basic, human-readable text format. Generated by most alignment algorithms. Consists of:

- **Header section (optional):** Lines start with @. Contains metadata: SAM format version (@HD), reference sequence dictionary (@SQ), read groups (@RG), program used (@PG), comments (@CO).
- **Alignment section:** Each line is an alignment record for a single read. Contains 11 mandatory fields, followed by optional fields.

Field Descriptions:

1. QNAME: Query template NAME
2. FLAG: bitwise FLAG
3. RNAME: Reference sequence NAME
4. POS: 1-based leftmost mapping POSition
5. MAPQ: MAPping Quality
6. CIGAR: CIGAR string
7. RNEXT: Ref. name of the mate/next read
8. PNEXT: Position of the mate/next read
9. TLEN: observed Template LENgth
10. SEQ: segment SEQuence
11. UAL: ASCII of Phred-scaled base QUALity+33

BAM:

Same format except that it is encoded in binary (faster to read) but not human legible.

CRAM:

Retains same info as SAM and is compressed in more efficient way.

Formats are *output* from aligners and assemblers.

BED Format:

Simple way to define basic sequence features to a sequence. One line per feature, each containing 3 - 12 columns of data plus optional track definition lines. Generally used for user defined sequence features as well as graphical representations of features.

- Chromosome Name
- Chromosome Start

- Chromosome End

Optional Fields Nine additional fields are optional for feature definition. If higher-numbered optional fields are used, all lower-numbered fields preceding them must also be populated.

Name Label to be displayed under the feature if enabled in the page configuration.

Score A numerical score ranging from 0 to 1000. The display style for scored data can be configured using track lines (see below).

Strand Defines the orientation of the feature: '+' (forward) or '-' (reverse).

thickStart The coordinate where the feature representation begins as a solid rectangle.

thickEnd The coordinate where the feature representation as a solid rectangle ends.

itemRgb An RGB color value (e.g., 0,0,255). This is applied only if a track line specifies `itemRgb="on"` (case-insensitive).

blockCount The number of sub-elements (e.g., exons) within the feature.

blockSize A comma-separated list of the sizes of these sub-elements.

blockStarts A comma-separated list of the start coordinates for each sub-element, relative to the feature's start coordinate.

Track Lines Track definition lines configure the display of features, such as grouping them into separate tracks. Track lines must precede the list of features they affect and consist of the word `track` followed by space-separated key=value pairs. Valid parameters for Ensembl include:

name A unique identifier for the track when parsing the file.

description A label displayed under the track in detailed views (e.g., "Region in Detail").

priority An integer determining the display order if multiple tracks are defined.

color Specified as RGB, hexadecimal, or an [X11 named color](#).

useScore A value from 1 to 4, dictating how scored data is displayed.
May require additional parameters:

- Tiling array
- Colour gradient (defaults to Yellow-Green-Blue with 20 grades). Optionally, custom colors (`cgColour1`, `cgColour2`, `cgColour3`) and the number of grades (`cgGrades`) can be specified.
- Histogram
- Wiggle plot

itemRgb If set to `on` (case-insensitive), the individual RGB values defined for each feature (in the `itemRgb` field) will be used.

BedGraph Format

The BedGraph format is designed for displaying moderate amounts of scored data and is based on the BED format with these key differences:

- The score is located in column 4 (instead of column 5 as in standard BED with score).
- Track lines are **compulsory** and must include `type=bedGraph`.

Optional track line parameters currently supported by Ensembl for BedGraph are:

name (as described above)

description (as described above)

priority (as described above)

graphType Specifies the display style, either `bar` or `points`.

June 20th / 2025

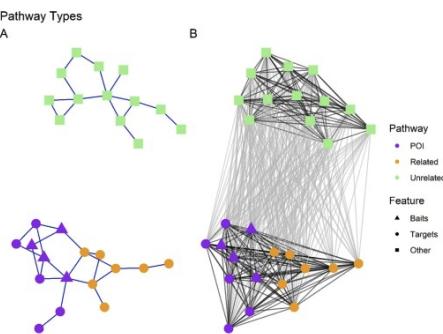
MASCARA: Coexpression analysis in data from designed experiments

13

MASCARA (Mixed-dAtatype State-space model for analyzing clinicAl tRajectories): Probabilistic framework designed to model longitudinal dynamics of multiple variables of mixed data types (e.g., continuous, categorical, count). The core of MASCARA is a State-Space Model (SSM), which posits that observed data are generated from an unobserved (latent) state that evolves over time.

Differential Expression Analysis (DE): Determines which features are expressed differently between 2 or more experimental conditions.

Coexpression Analysis (CoE): Aims at discovering features that are part of an already partially characterized pathway of interest (POI). the expression profiles of known features of POI (baits) are used to detect novel pathways members with similar expresssuin or accumulation patterns (targets).



Mathematical Framework

SSM consists of 2 primary components: state equation and observation equation

¹³ Fred T.G. White, Anna Heintz-Buschart, Lemeng Dong, Harro J. Bouwmeester, Johan A. Westerhuis, and Age K. Smilde. Mascara: Coexpression analysis in data from designed experiments. *Computational and Structural Biotechnology Reports*, 2: 100052, 2025. ISSN 2950-3639. doi: <https://doi.org/10.1016/j.csbr.2025.100052>. URL <https://www.sciencedirect.com/science/article/pii/S2950363925000237>

State Equation (Transition Model)

The evolution of the system is captured by a low-dimensional latent state vector $\mathbf{z}_t \in \mathbb{R}^L$ at each time point t . MASCARA models this evolution as a linear Gaussian process, such as a random walk.

$$\mathbf{z}_t | \mathbf{z}_{t-1} \sim \mathcal{N}(\mathbf{A}\mathbf{z}_{t-1}, \mathbf{Q}) \quad (20)$$

where \mathbf{A} is the transition matrix and \mathbf{Q} is the process noise covariance matrix. For a simple random walk, \mathbf{A} is the identity matrix.

Observation Equation (Emission Model)

The link between the latent state \mathbf{z}_t and the observed variable $y_{j,t}$ is modeled using a generalized linear model framework [2]. For categorical variables, which are most relevant to the ClinVar analysis, a multinomial logistic regression model is used.

Probability of observing category k for j -th variable at time t :

$$P(y_{j,t} = k | \mathbf{z}_t) = \text{softmax}(\mathbf{W}_j \mathbf{z}_t + \mathbf{b}_j)_k \quad (21)$$

For a vector \mathbf{x} :

$$\text{softmax}(\mathbf{x})_k = \frac{\exp(\mathbf{x}_k)}{\sum_{l=1}^K \exp(\mathbf{x}_l)} \quad (22)$$

\mathbf{W}_j and \mathbf{b}_j learned during model fitting

Model Inference

The model parameters and latent states are estimated using **Stochastic Variational Inference (SVI)** makes method scalable to large datasets. PyTorch SVI framework.

June 25th / 2025

Simple Conceptual Pipeline

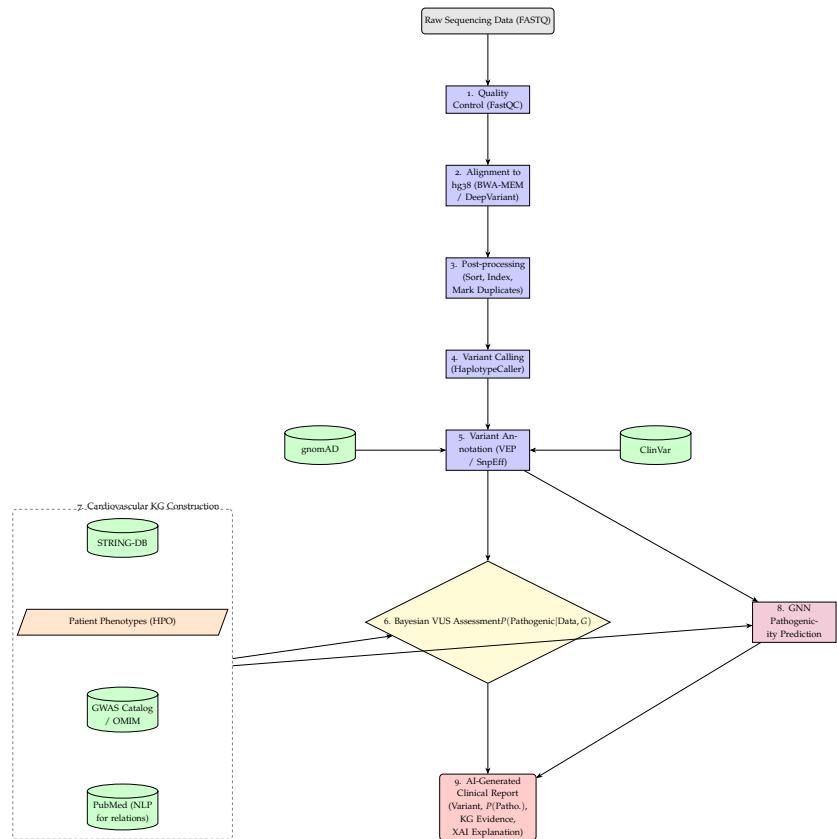


Figure 3: Conceptual pipeline example.

Dataflow Programming

Paradigm that models a program as a directed graph of data flowing between operations. Specify independent tasks and the data de-

pendencies between them, The program executes as data becomes available.

Components

- **Processes:** Operations/ Nodes in graph. Self-contained independent script that performs x computational task. (e.g., command-line tool, Python / R scripts,...,)
- **Channels:** Data Pipes / Edges. Asynchronous queues that connect processes. (sort of like a NN block or something like that).

Flow Data

Driven entirely by the availability of data channels.

- **Initiation:** When ≥ 1 initial channels are populated. Channel Factory “Channel.somepath”
- **Process Execution**
- **Implicit Parallelism**

e.g., “samtools” pipeline

```
Input Channel (e.g., a BAM file)
  |
  v
[ samtools_sort ]
  |
  +-----+-----+
  |       |       |
  v       v       v
[ samtools_idx ] [ samtools_stats ] [ samtools_flagstat ]
```

Encoding word order in complex embeddings

¹⁴

¹⁴ NYU Langone Health. Next-generation sequencing analysis resources, 2023. URL <https://learn.genome.bio.nyu.edu>. Accessed on 2025-05-30

June 27th / 2025

AlphaGenome: advancing regulatory variant effect prediction with a unified DNA sequence model

15

Challenge: Interpretation of non-coding genetic variants - 98% of human genetic variation and difficult to understand due to diverse / subtle effects on gene regulation.

Trade-off between *length* and *prediction resolution* due to computational limitations

- **High Resolution, short sequence length:** Some models (e.g., BP-Net, SpliceAI) achieve high nucleotide-level predictive resolution at the **cost** of being restricted to short input DNA sequences 10kb or less. **CANNOT** capture long-range genomic interactions and miss influence of important distal regulatory elements that lie outside limited input window
- **Long Sequence Length, low resolution:** Models like **Enformer** and **Borzoi** can process long DNA seq. ($\approx 200\text{kb}$ to $\approx 500\text{kb}$), which allows them to capture a broader regulatory context. The trade-off is resolution of predictions must be reduced (e.g., group output into bins of 32 or 128 bps). May miss blur fine features
- **Multimodal Predictions:** Prediction Ks genomic features across 11 different modalities (RNA-seq, CAGE), chromatin accessibility (ATAC-seq, DNase-seq), histone modifications, transcription factor binding, chromatin contact maps (Hi-C) and detailed splicing patterns (splice, usage, junctions)
- **Splicing:** First model to predict splice sites, splice site usage, and splice junctions counts simultaneously → How do variants disrupt splicing?
- Foundation of AlphaGenome is vast / diverse collection of public functional genomics data from human (hg38 reference genome)

¹⁵ Žiga Avsec, Natasha Latysheva, Jun Cheng, Guido Novati, Kyle R. Taylor, Tom Ward, Clare Bycroft, Lauren Nicolaisen, Eirini Arvaniti, Joshua Pan, Raina Thomas, Vincent Dutordoir, Matteo Perino, Soham De, Alexander Karollus, Adam Gayoso, Toby Sargeant, Anne Mottram, Lai Hong Wong, Pavol Drotár, Adam Kosiorek, Andrew Senior, Richard Tanburn, Taylor Applebaum, Souradeep Basu, Demis Hassabis, and Pushmeet Kohli. AlphaGenome: advancing regulatory variant effect prediction with a unified DNA sequence model. *bioRxiv*, 2024. doi: 10.1101/2024.03.01.583020. URL <https://www.biorxiv.org/content/10.1101/2024.03.01.583020v1>

and mouse (mm10 reference genome) sources. Process was standardized to ensure high quality and consistency.

- **Data Sources and Standardization**

- Aggregation from ENCODE, FANTOM5, GTEx:
 - * RNA-seq, CAGE, PRO-cap, DNase-seq, ATAC-seq, ChIP-seq, chromatin contact maps
- Standardize metadata for ALL samples:
 - * UBERON for tissues, EFO/CLO for cell lines, CL for primary cell types
- QC step for all ENCODE = application of ENCODE's audit system

Processing of Specific Data Types

(More specific)

RNA-seq

- Data from both ENCODE and GTEx (via RECOUNT₃) were used.
- To ensure comparability, all tracks were normalized to Reads Per Million (RPM) and then further rescaled to a common factor of 100 million total reads.
- Tracks were then grouped by biological context (ontology term, assay type, etc.) and averaged to create final representative tracks for model training.

CAGE and PRO-cap

- CAGE data from FANTOM5 and PRO-cap data from ENCODE were processed similarly to RNA-seq, with individual tracks normalized to a total of 100 million reads before any averaging.

DNase-seq and ATAC-seq (Chromatin Accessibility)

- Instead of using pre-processed signal files, raw alignment (BAM) files were used to preserve base-resolution information.
- These were converted to bigWig files that record the counts of enzyme cut sites at each base.
- This approach allows for the correction of enzyme cut bias by applying appropriate read shifts (+4 / -4 for ATAC-seq and 0 / +1 for DNase-seq).
- The resulting tracks were averaged within ontology groups and normalized to 100 million insertions per track.

ChIP-seq (TF and Histone)

- Only "fold change over control" signal files were selected.
- After stringent QC and hierarchical filtering to select the most representative experiments, the fold-change signals were averaged within each biological context group without further normalization.

Splicing Data (Junctions, Sites, Usage)

- RNA-seq reads were realigned using STAR to specifically detect and quantify splice junctions.
- **Splice Junctions:** Raw junction counts were subjected to a stringent quality filtering pipeline to ensure high confidence. The retained counts were then normalized and scaled for model training.
- **Splice Site Usage (SSU):** SSU was calculated for each potential splice site. The quantification was performed considering all reads spanning the splice sites regardless of the strand. The formula used is:

$$\text{SSU} = \frac{\text{\#reads using the splice site}}{\text{\#reads using the splice site} + \text{\#reads supporting skipping of the splice site}} \quad (23)$$

- **Splice Site Classification:** The set of all unique donor and acceptor sites from the filtered junction data was used to define the training examples for a 5-class classification task (Donor+, Acceptor+, Donor-, Acceptor-, or Not a splice site).

Contact Maps

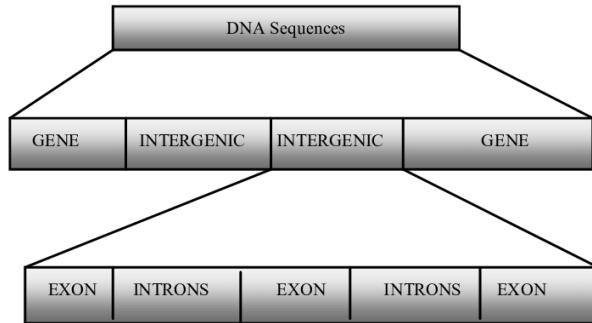
- Hi-C and Micro-C datasets were sourced from the 4D Nucleome portal.
- The maps were processed following the Orca protocol, which includes matrix balancing and adaptive coarse-graining.
- A distance-based normalization was applied to compute a log-fold change over the average contact value for each genomic distance. The normalized contact map values $y_{i,j}$ were computed as:

$$y_{i,j} = \log \left(\frac{x_{i,j} + \epsilon}{\text{mean}[|i - j|] + \epsilon} \right) \quad (24)$$

where $x_{i,j}$ is the coarse-grained count, ϵ is a numerical relaxation constant, and $\text{mean}[|i - j|]$ is the average coarse-grained contact value for the pairwise genomic distance.

Splicing

Process that generates mature mRNA sequences. It functions by removing non-coding regions (**introns**) and joining remaining coding regions (**exons**)



Genetic variants may disrupt splicing which alters final mRNA → aberrant protein

- **Splice Sites:** Specific nucleotides in DNA seq. that signal where splicing machinery should *cut and join* RNA. They mark boundaries between introns and exons:
 - **Splice Donors:** beginning of intron
 - **Splice acceptor:** End of intron

Probability any given nucleotide will function as donor / acceptor can be modeled and such sites are associated with *recognizable sequence motifs*. **αGenome** predicts classification of these sites on **positive** and **negative** DNA strands (wtf is this?)

Splice Junction

Specific connection point where 2 exons are ligated (joined) together after having intron excised. Prediction of which introns are removed = **Splice Junction Prediction**

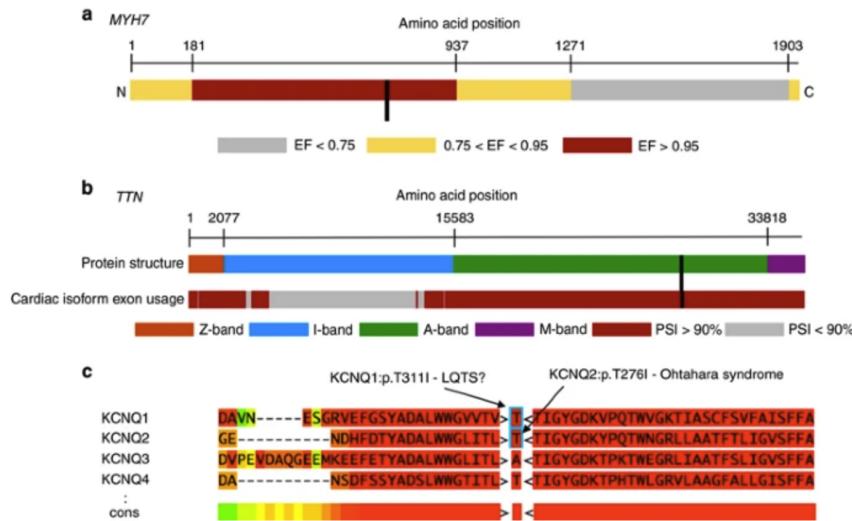
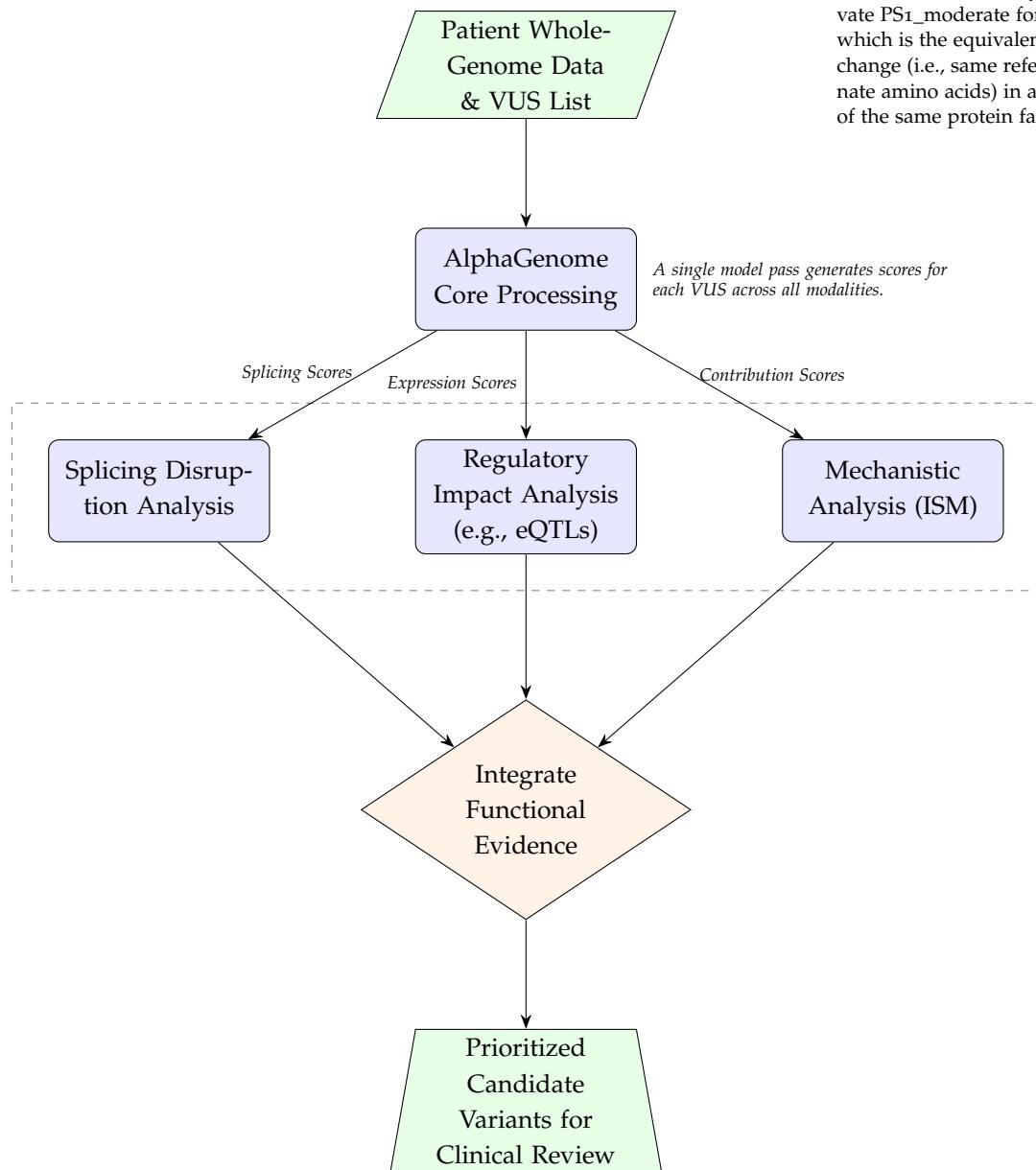


Figure 4: (a) Missense variants within a subportion of *MYH7*, when identified in an HCM patient, have a 97% prior probability of being pathogenic (etiological fraction; EF = 0.97). We activate PM1 for missense variants in this region. Here we use *MYH7:c.2221G>T* as an example (black bar). (b) Truncating variants in *TTN* are only known to cause DCM when found in exons constitutively expressed in the heart (proportion spliced in >0.9). We activate PVS1_strong for these variants. Here we use *TTN:c.8664_1delC* as an example (black bar). (c) Variants that have been identified as pathogenic in paralogous genes may identify residues that are intolerant to variation. We have created two modified rules, PS1_moderate and PM5_supporting, to incorporate this evidence. Here we use *KCNQ1:p.T311I* as an example. *KCNQ2:p.T276I* is associated with Ohtahara syndrome. We activate PS1_moderate for *KCNQ1:p.T311I*, which is the equivalent missense change (i.e., same reference and alternate amino acids) in a different member of the same protein family.



June 28th / 2025

RECAP:

- Pathogenic (P): Sufficient evidence to be considered disease-causing
- Likely Pathogenic (LP): Evidence strongly suggests pathogenicity, but is not definitive
- Variant of Uncertain Significance (VUS): Insufficient or conflicting evidence to classify as either pathogenic or benign
- Likely Benign (LB): Evidence strongly suggests a benign (non-disease-causing) impact
- Benign (B): Sufficient evidence to be considered not disease-causing

Such framework is fundamentally probabilistic. 'Likely' categories > 90% certainty. Inherent probabilistic must be captured by predictive model!

Such guidelines have shown to reduce reclassification rates.

How Evidence Drives Reclassification?

Reclassification rates ranges from 3.6% to 58.6% (really retarded range) underscoring the *volatility* of initial classifications and importance of a continuous influx of new data.

Accumulation of information from multiple sources drives classification:

- **Functional Studies** such as multiplex assays of variant effects (MAVEs) measure functional impact of variant on protein function (time consuming and resource intensive)
- **Computational / In Silico Data** AlphaMissense provides novel in silico evidence for reclassification (analysis of evolutionary conservation and predicted structural impact)

- **Population Data** gnomAD provides frequency of variants in general population. e.g., variant common in healthy subjects is less likely to be pathogenic for a rare disease
- **Segregation Data** family studies / segregation studies to track whether x variant co-occurs with a disease across multiple family members
- **Literature and case reports** linking variants to phenotype

VUS

VUS imply a lack of compelling evidence to confidently place a variant under a category → Bottleneck in clinical genetics. A model that fails to predict a VUS-to-Pathogenic transition has made a more clinically significant error than one that misses an LP-to-P transition → A model's loss function could be customized to more **heavily penalize errors** optimizing model to align with clinical utility.

Categorical Time-Series Problem

The reclassification of variants can be seen as a **Time-ordered sequence of categorical labels**. e.g., Classification of variant x at time t is **dependent** on classification at $t - 1 +$ cumulative evidence at $t \rightarrow$ Temporal Dependency. Fundamental problem: Ground truth is *non-stationary*. Whatever constitutes a pathogenic variant is refined over time, so statistical properties of data-generating process are *changing*.

Concept Drift: Phenomenon where statistical relationships between input data and target variables change over time, usually in dynamic environments. *Patterns a machine learning model learned from historical data may no longer be valid as the real-world environment evolves.*
So, how to solve such problem?



Framing a Predictive Task as Concept Drift

Simple time-series forecasting models assume that underlying patterns in data are stable over time. Variant classification violates such assumption! since its intrinsic properties and designed classification evolves as knowledge grows. A

A **stationary process** is one in which μ and σ^2 remain constant over time. Therefore, variant classification is an inherently non-stationary process. Concept drift = non-stationary process but in supervised learning context → *Change over time in the statistical relationship between input X and target variable Y, $P(y|X)$* . Intrinsic features of a variant (e.g., genomic location, nucleotide change...) are static BUT interpretation of such features (concept of pathogenicity) evolves. So NOT the same as **Data Drift** where $P(X)$ changes.

Drift Patterns

- **Abrupt Drift:** rapid and sudden change in data distribution. It could be triggered by a singular, high impact event, e.g., publication of 2015 ACMG/AMP guidelines
- **Gradual Shift:** Slow, incremental change over a long period. *Likely observed in variant classification.* Evidence accumulates over time, slowly shifting balance of evidence for or against pathogenicity
- **Recurring Drift:** Cyclical changes. Less common in variant reclassification.

The timestamp of reclassification serves as proxy for state of scientific knowledge and classification standards in place at t . Thus, time-related features (e.g., time elapsed since last classification, binary indicators for guideline changes) MUST be incorporated into model's input.

A **Stream Learning** framework *MUST* be so it can adapt and evolve as new data becomes available.

Strategies for Model Updating

- **Periodic Retraining** Retained from scratch at regular intervals (every x months) using data up to t
- **Triggered Retraining** COntinuously monitor model's performance on new incoming data. An algorithm such as **Drift Detection Method (DDM)** or **LSTM-based detector (LSTMDD)** is used to identify a statistically significant drop in performance. When drift is detected → Trigger MODEL RETRAINING. NO unnecessary monthly/annual updates while '*concept*' = *stable*

THEREFORE, model must be able to learn historical patterns AND patterns of *how those patterns themselves change!*. An LSTM will capture temporal dependencies BUT it MUST be embedded within a larger online framework that manages its lifecycle of training, evaluation, and updating.

Temporal Long Short-Term Memory (LSTM)

Before LSTM learns from temporal sequences of variant reclassification:

1. encode categorical classification labels into numerical format
2. transform continuous-time series data into discrete input-output windows suitable for supervised learning

Entity Embedding

From NLP (e.g., Word2Vec), it maps each category to a dense, low dimensional vector of continuous values (for each category within a categorical feature). Vectors are not predefined but learned by NN. An embedding layer is added to the model which learns to place categories that are semantically similar closer to each other in the multi-dimensional embedding space → Model can autonomously discover that LP and P should have similar vector representations, while B should be distant from them. Entity embeddings thus provide dense and **contextually** aware representation. Also good for sparse or complex categorical data.

Embedding Layer: Core component of entity embeddings. It is pretty much a lookup table which is represented as a weight matrix.

Embedding Matrix

Let C be a categorical feature with a vocabulary of size V . This means the feature has V unique categories (e.g., if the feature is "Gene Symbol", V would be the total number of unique gene symbols in the dataset).

Let D be the desired dimensionality of embedding vectors. This is a hyperparameter (e.g., $D = 10, D = 50$)

The embedding layer maintains a single weight matrix, **embedding matrix**, E .

$$E \in \mathbb{R}^{V \times D} \quad (25)$$

Each row e_i of matrix E is the D -dimensional embedding vector for the i -th category in the vocabulary.

$$E = \begin{bmatrix} \cdots & e_1 & \cdots \\ \cdots & e_2 & \cdots \\ \vdots & \vdots & \vdots \\ \cdots & e_V & \cdots \end{bmatrix} = \begin{bmatrix} e_{11} & e_{12} & \dots & e_{1D} \\ e_{21} & e_{22} & \dots & e_{2D} \\ \vdots & \vdots & \ddots & \vdots \\ e_{V1} & e_{V2} & \dots & e_{VD} \end{bmatrix}$$

Initially, matrix is filled with small random values. The goal of training is to learn optimal values for all entries in matrix

Embedding Lookup Operation

The input to an embedding layer is typically an integer representing a specific category (e.g., category "Gene Symbol: BRCA1" might be mapped to integer index 15). The layer performs a simple **lookup operation** to retrieve the corresponding embedding vector.

Lookup can be represented as a matrix multiplication with a one-hot encoded vector. Let the input category be represented by integer k , where $1 \leq k \leq V$. Then represent this input as a one-hot vector c_k of size V , where k -th element is 1 and all other elements are 0.

$$c_k = [0, 0, \dots, 1, \dots, 0]^T \quad (26)$$

The output of the embedding layer, which is the embedding vector for category k , is then given by:

$$v_k = E^T \cdot c_k \quad (27)$$

This matrix multiplication selects the k -th column of E^T , which is the k -th row of E . NNs don't go through all this BS with matrices, this is just for simplification.

Learning Embeddings

Values of embedding matrix E are learned (backpropagation). The network is trained on a supervised task (e.g., predicting variant pathogenicity).

Objective Function (Loss Function)

Let the neural network be denoted by a function f , which takes the learned embeddings and other features as input to make a prediction \hat{y} . The network learns by minimizing a loss $L(\hat{y}, y)$

For a classification task, a common loss function is the **Binary Cross-Entropy Loss**:

$$L(\hat{y}, y) = -\frac{1}{N} \sum_{i=1}^N [y_i \log(\hat{y}_i) + (1 - y_i) \log(1 - \hat{y}_i)] \quad (28)$$

$\hat{y}_i = f(v_{k_i}, \dots)$ is output for i -th training, which includes embedding vector v_{k_i} for category k

Updating Embedding Matrix (Backpropagation)

Use of optimization algorithm (e.g., Stochastic Gradient Descent (SGD), Adam) to update all weightsn and embedding matrix E .

For a single training example, the gradient of the loss L with respect to the embedding matrix E is calculated. An important property of the embedding lookup operation is that the gradient is **sparse**. For an input category k , only k -th row of the embedding matrix, e_k , contributed to the final output. Therefore, the gradient will be non-zero only for that specific row.

$$\frac{\partial L}{\partial e_j} = 0 \quad \text{for all } j \neq k \quad (29)$$

The non-zero gradient is for the vector e_k that was used in the forward pass:

$$\frac{\partial L}{\partial e_k} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial e_k} \quad (30)$$

Update e.g., standard SGD with learning rate η :

$$e_k^{(\text{new})} = e_k^{(\text{old})} - \eta \frac{\partial L}{\partial e_k} \quad (31)$$

Over many training iterations, each row of the embedding matrix is adjusted based on how it contributes to the overall prediction error for the categories it represents. This process causes vectors for categories that appear in similar contexts (with respect to the prediction target) to be pushed closer together in the embedding space.

Determining the optimal window size ' n_{past} ' is a critical hyperparameter tuning task.

- Most variant classifications occur within ≈ 2 years of initial classification. If data is recorded quarterly then window size of 8 (2 years) is a logical initial value for test
- Treat window size as hyperparameter and evaluate performance *dedicated validation set*

Validation for Temporal Data

For data with temporal dependencies, such as variant reclassification histories standard validation techniques are inappropriate guaranteed to produce misleading and way too optimistic results

Let a time-series dataset be an ordered sequence of observations $D = \{(x_1, y_1), (x_2, y_2), \dots, (x_T, y_T)\}$, where T is total number of observations.

The process creates a series of k splits, where each split consists of a training set ($D_{\text{train}}^{(i)}$) and a validation set ($D_{\text{val}}^{(i)}$).

Splitting

Let n_0 be the size of the initial training set, and let h be the size of the validation set (the forecast horizon). The splits are generated iteratively.

For each split $i = 1, 2, \dots, k$:

- **Training Set ($D_{\text{train}}^{(i)}$):** The training set includes all data from the beginning up to a certain point in time. It expands with each iteration.

$$D_{\text{train}}^{(i)} = \{(x_t, y_t) \mid t = 1, \dots, n_0 + (i-1)h\} \quad (32)$$

- **Validation Set ($D_{\text{val}}^{(i)}$):** The validation set consists of the next h observations immediately following the training set.

$$D_{\text{val}}^{(i)} = \{(x_t, y_t) \mid t = (n_0 + (i-1)h) + 1, \dots, n_0 + ih\} \quad (33)$$

The constraint $n_0 + k \cdot h \leq T$ must be satisfied. This ensures that the validation set for the final split does not exceed the total number of observations.

Model Evaluation

A model, M , is trained on each training set $D_{\text{train}}^{(i)}$ to produce a trained model M_i . The performance of each model M_i is then evaluated on its corresponding validation set $D_{\text{val}}^{(i)}$ using a chosen error metric, \mathcal{L} (e.g., Mean Squared Error, Mean Absolute Error). Let this error be Err_i .

$$\text{Err}_i = \mathcal{L}(M_i, D_{\text{val}}^{(i)}) \quad (34)$$

The overall performance estimate for the model M is the average of the errors across all k splits.

$$\text{CV}_{\text{Error}} = \frac{1}{k} \sum_{i=1}^k \text{Err}_i \quad (35)$$

This procedure ensures that at no point is the model trained on data that occurred chronologically after the data it is being asked to predict, thus preserving the temporal integrity of the dataset.

Nested Cross-Validation is probably better than the above strategy!

Probabilistic Forecasting

Goal is to predict single most likely next classification for x variant. Model's uncertainty is paramount. e.g., A prediction of "Pathogenic" with 51% confidence carries a vastly different clinical implication than a prediction with 99% confidence. Therefore, the task should be framed as **probabilistic classification** where model outputs a full probability distribution over all possible future classes

Architecture must be designed as follows:

- **Output Layer:** Final Layer of LSTM should be a *Dense (fully connected)* layer
- **Activation Function:** Dense layer must use **softmax activation function**. in: logits, out: probability vector
- **Number Neurons:** Must equal to number of classes in classification system (B, LB, VUS, LP, P).
e.g., Input sequence = vector [0.05, 0.10, 0.10, 0.60, 0.15], representing a 5% chance of being Benign, 10% Likely Benign, 10% VUS, 60% Likely Pathogenic, and 15% Pathogenic
 - **Attention Mechanism:** Attention layer to allow LSTM to dynamically assign different weights of important to different time steps in the input sequence when making a prediction.
 - **Stacked LSTM:** Multiple layers can be stacked on top of each other to allow model to learn a hierarchical representation of the temporal data. e.g., first layer may learn simple short-term patterns, while deeper layers could combine these to learn more complex, long-term abstract patterns in the reclassification sequence

Bibliography

Žiga Avsec, Natasha Latysheva, Jun Cheng, Guido Novati, Kyle R. Taylor, Tom Ward, Clare Bycroft, Lauren Nicolaisen, Eirini Arvaniti, Joshua Pan, Raina Thomas, Vincent Dutordoir, Matteo Perino, Soham De, Alexander Karollus, Adam Gayoso, Toby Sargeant, Anne Mottram, Lai Hong Wong, Pavol Drotár, Adam Kosiorek, Andrew Senior, Richard Tanburn, Taylor Applebaum, Souradeep Basu, Demis Hassabis, and Pushmeet Kohli. AlphaGenome: advancing regulatory variant effect prediction with a unified DNA sequence model. *bioRxiv*, 2024. DOI: 10.1101/2024.03.01.583020. URL <https://www.biorxiv.org/content/10.1101/2024.03.01.583020v1>.

Jun Cheng, Guido Novati, Joshua Pan, Clare Bycroft, Akvilė Žemgulytė, Taylor Applebaum, Alexander Pritzel, Lai Hong Wong, Michal Zielinski, Tobias Sargeant, Rosalia G. Schneider, Andrew W. Senior, John Jumper, Demis Hassabis, Pushmeet Kohli, and Žiga Avsec. Accurate proteome-wide missense variant effect prediction with alphamissense. *Science*, 381(6664):eadg7492, 2023. DOI: 10.1126/science.adg7492. URL <https://www.science.org/doi/abs/10.1126/science.adg7492>.

Yanrong Ji, Zhihan Zhou, Han Liu, and Ramana V Davuluri. Dnabert: pre-trained bidirectional encoder representations from transformers model for dna-language in genome. *Bioinformatics*, 37(15):2112–2120, 02 2021. ISSN 1367-4803. DOI: 10.1093/bioinformatics/btab083. URL <https://doi.org/10.1093/bioinformatics/btab083>.

J. Jumper, R. Evans, A. Pritzel, and et al. Highly accurate protein structure prediction with alphafold. *Nature*, 596(7873), 2021.

Zhen Liu, Yifan Gu, and Xiaoyang Huang. Deep learning-based ranking method for subgroup and predictive biomarker identification in patients. *Communications Medicine*, 5:221, 2025. DOI: 10.1038/s43856-025-00946-z. URL <https://doi.org/10.1038/s43856-025-00946-z>.

L Meng, R Attali, T Talmy, Y Regev, N Mizrahi, P Smirin-Yosef, L Vossaert, C Taborda, M Santana, I Machol, R Xiao, H Dai, C Eng, F Xia, and S Tzur. Evaluation of an automated genome interpretation model for rare disease routinely used in a clinical genetic laboratory. *Genet Med*, 25(6):100830, 2023. DOI: [10.1016/j.gim.2023.100830](https://doi.org/10.1016/j.gim.2023.100830).

V Nauffal, P Di Achille, M D R Klarqvist, and et al. Genetics of myocardial interstitial fibrosis in the human heart and association with disease. *Nature Genetics*, 55:777–786, 2023. DOI: [10.1038/s41588-023-01371-5](https://doi.org/10.1038/s41588-023-01371-5).

P. C. Ng and S. Henikoff. Sift: predicting amino acid changes that affect protein function. *Nucleic Acids Research*, 31(13):3812–3814, 2003. DOI: [10.1093/nar/gkg509](https://doi.org/10.1093/nar/gkg509).

NYU Langone Health. Next-generation sequencing analysis resources, 2023. URL <https://learn.gencore.bio.nyu.edu>. Accessed on 2025-05-30.

S Richards, N Aziz, S Bale, D Bick, S Das, J Gastier-Foster, W W Grody, M Hegde, E Lyon, E Spector, K Voelkerding, H L Rehm, and ACMG Laboratory Quality Assurance Committee. Standards and guidelines for the interpretation of sequence variants: a joint consensus recommendation of the american college of medical genetics and genomics and the association for molecular pathology. *Genetics in Medicine*, 17(5):405–424, 2015. DOI: [10.1038/gim.2015.30](https://doi.org/10.1038/gim.2015.30).

Wuwei Tan and Yang Shen. Multimodal learning of noncoding variant effects using genome sequence and chromatin structure. *Bioinformatics*, 39(9):btad541, 09 2023. ISSN 1367-4811. DOI: [10.1093/bioinformatics/btad541](https://doi.org/10.1093/bioinformatics/btad541). URL <https://doi.org/10.1093/bioinformatics/btad541>.

Fred T.G. White, Anna Heintz-Buschart, Lemeng Dong, Harro J. Bouwmeester, Johan A. Westerhuis, and Age K. Smilde. Mascara: Coexpression analysis in data from designed experiments. *Computational and Structural Biotechnology Reports*, 2:100052, 2025. ISSN 2950-3639. DOI: <https://doi.org/10.1016/j.csbr.2025.100052>. URL <https://www.sciencedirect.com/science/article/pii/S2950363925000237>.