

BRYAN L.A.

# LITERATURE NOTES BIOINFORMATICS

UNIVERSITY OF AMSTERDAM  $\cap$  VRIJE UNIVERSITEIT

No copyright © 2025 Bryan L.A.

UNIVERSITY OF AMSTERDAM  $\cap$  VRIJE UNIVERSITEIT

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

*Last updated June 2025*

# *Contents*

*June 11<sup>th</sup> / 2025*      5

*June 12<sup>th</sup> / 2025*      13

*June 13<sup>th</sup> / 2025*      21

*Bibliography*      23



June 11<sup>th</sup> / 2025

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

**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

---

<sup>1</sup> 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

*Methods:*

- Supervised learning approach, trained on dataset of 10<sup>3</sup>'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
  - **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 PATIENTS <sup>2</sup>

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.

<sup>2</sup> 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>

## 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 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  $p \gg \gg$ , 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:

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

- $\alpha$  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,  $\phi$  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:

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

- $\theta$  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:**  $\theta$  (including  $\alpha$  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  $\rightarrow$  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  $\rightarrow$  ranking by  $\hat{\tau}(x)$  separates patients who *truly* benefit from those who do not.

### *Mathematical Distinction*

**Prognostic Model (Covariates  $\rightarrow$  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  $\beta_w$  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  $\rightarrow$  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 12<sup>th</sup> / 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<sup>3 4</sup>

### *$\alpha$ -Missense*

By Google DeepMind. It predicts pathogenicity of missense variants. Architecture is inspired by *Evoformer* block used in  $\alpha$ -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  $S_{target} \in \{0, 1\}^{L \times 20}$ , 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 refines 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)

<sup>3</sup> 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>

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

### ***The Block:***

*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  $\rightarrow$  **Iterative refinement**(reasoning cycle). Model uses *current* pair representation  $P$  to *bias/guide* attention mechanism. Thus instead of e.g.,  $\text{Scores} = QK^T$ , it becomes  $\text{Scores} = QK^T + \text{Bias}_{ij}$ . A large  $\text{Bias}_{ij}$  = 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

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

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

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**.

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

$\alpha$  is final pathogenicity score



## 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 GENOME<sup>6</sup>

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.

<sup>6</sup> 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>

## 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,  $\gamma_i$ :

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

Assumes variant effects  $\beta_j$  for each category  $k$  are random variables drawn from distribution with a category-specific variance component,  $\tau_k$ :

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

$H_0$  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  $\chi^2$  distributions from which p-value can be derived.

---

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



*June 13<sup>th</sup> / 2025*



# Bibliography

- 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.
- 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.