

Linkage Disequilibrium-Informed Deep Learning Framework to Identify Genetic Loci for Alzheimer's Disease Using Whole Genome Sequencing Data

Taeho Jo¹, Paula Bice¹, Kwangsik Nho^{1,2*}, Andrew J. Saykin^{1,3*} and the Alzheimer's Disease Sequencing Project

1. Indiana Alzheimer Disease Research Center and Center for Neuroimaging, Department of Radiology and Imaging Sciences, Indiana University School of Medicine, Indianapolis, IN 46202, USA
2. Center for Computational Biology and Bioinformatics, Indiana University School of Medicine, Indianapolis, IN 46202, USA
3. Department of Medical and Molecular Genetics, Indiana University School of Medicine, Indianapolis, IN 46202, USA

Abstract

The exponential growth of genomic datasets necessitates advanced analytical tools to effectively identify genetic loci from large-scale high throughput sequencing data. This study presents Deep-Block, a multi-stage deep learning framework that incorporates biological knowledge into its AI architecture to identify genetic regions as significantly associated with Alzheimer's disease (AD). The framework employs a three-stage approach: (1) genome segmentation based on linkage disequilibrium (LD) patterns, (2) selection of relevant LD blocks using sparse attention mechanisms, and (3) application of TabNet and Random Forest algorithms to quantify single nucleotide polymorphism (SNP) feature importance, thereby identifying genetic factors contributing to AD risk. The Deep-Block was applied to a large-scale whole genome sequencing (WGS) dataset from the Alzheimer's Disease Sequencing Project (ADSP), comprising 7,416 non-Hispanic white participants (3,150 cognitively normal older adults (CN), 4,266 AD). First, 30,218 LD blocks were identified and then ranked based on their relevance with Alzheimer's disease. Subsequently, the Deep-Block identified novel SNPs within the top 1,500 LD blocks and confirmed previously known variants, including *APOE* rs429358 and rs769449. The results were cross-validated against established AD-associated loci from the European Alzheimer's and Dementia Biobank (EADB) and the GWAS catalog. The Deep-Block framework effectively processes large-scale high throughput sequencing data while preserving interactions between SNPs in performing the dimensionality reduction, which can potentially introduce bias or lead to information loss. The Deep-Block approach identified both known and novel genetic variation, enhancing our understanding of the genetic architecture of and demonstrating the framework's potential for application in large-scale sequencing studies.

Introduction

The advancement of deep learning in artificial intelligence has introduced new frameworks for analyzing complex genetic inheritance patterns, enhancing the interpretation of genomic data. (Avsec et al., 2021; Eraslan et al., 2019; Novakovsky et al., 2023; Zhou & Troyanskaya, 2015). For complex diseases such as Alzheimer's disease (AD), there is a critical need for advanced analytic tools provided by Artificial Intelligence (AI) to decipher the complexities of human genetic makeup (Berson et al., 2023; Eraslan et al., 2019; Shigemizu et al., 2023). The complexity of genomic studies necessitates innovative and adaptive approaches that transcend traditional machine learning techniques to analyze and elucidate these intricate genetic interactions (Bettencourt et al.; Karczewski & Snyder, 2018). The high dimensionality and large sample sizes characteristic of genetic data in AD research underscore the necessity for methods capable of navigating the complex landscape (Jo et al., 2023; Jo et al., 2022; Konietschke et al., 2021). While several machine learning-based dimensionality reduction methods have been proposed, they have encountered challenges such as loss of phenotypic association information during the reduction process, reproducibility issues, and data-dependent inconsistencies in results (Fujiwara et al., 2020; Shetta & Niranjana, 2020; Vogelstein et al., 2021).

Here, we present Deep-Block, a deep learning framework designed to address the complexities of genomic sequencing data through targeted analysis of whole genome sequencing (WGS) data. Deep-Block employs a linkage disequilibrium (LD) block-based approach to systematically identify significant genetic regions, aiming to preserve vital phenotypic associations and minimize the loss of genetic information crucial for understanding disease phenotypes. A key feature of Deep-Block is to efficiently handle missing data, a common challenge in large-scale genetic studies. The framework incorporates advanced machine learning and genomic imputation techniques (An et al., 2023; Rubinacci et al., 2020; Shishegar et al., 2021) to ensure a comprehensive dataset without any missing values for analysis. Furthermore, the integration of the TabNet model (Arik & Pfister, 2021; Vaswani et al., 2017), an attention-based neural network, enhances the process by providing a detailed assessment of feature importance within the genetic data, thus enriching the analysis. The calculation of phenotype influence scores (PIS) offers additional insights into the genetic basis of the disease, informing future research directions.

Application of Deep-Block to a large-scale WGS dataset from the Alzheimer's Disease Sequencing Project (ADSP) Release 3, comprising 7,416 non-Hispanic white participants, demonstrated its capacity to effectively manage complex genomic data and identify single nucleotide polymorphisms (SNPs) as associated with AD. The Deep-Block framework identified AD-associated genetic loci, including well-

known AD SNPs such as *APOE* rs429358 and rs769449 and novel single nucleotide polymorphisms (SNPs) not previously reported in AD genetic association studies, particularly within the top-performing LD blocks.

Methods

Data Collection and Quality Control

The ADSP participants have WGS data sequenced using multiple platforms, including IlluminaHiSeq2000 and IlluminaHiSeqXTen. This release (R3) includes 16,906 whole-genome sequences (WGS), processed and curated as part of the project. The release contains CRAMs, gVCFs, and quality-controlled project-level VCFs (pVCFs) for autosomal biallelic single nucleotide variants (SNVs) and indels, along with structural variant (SV) calls generated by Manta, Smoove, and Strelka variant callers. The WGS data were called by the Genome Center for Alzheimer's Disease (GCAD) using VCPA 1.1, a functionally equivalent CCDG/TOPMed pipeline. WGS data underwent comprehensive quality control (QC) procedures, including SNP call rates > 95%, Hardy-Weinberg equilibrium P values < 1×10^{-6} , minor allele frequencies (MAF) > 1%, absence of sex mismatches, and sample call rates > 95%. To mitigate false associations due to population stratification, the study analyzed genome-wide genotyping data from 7,416 non-Hispanic White (NHW) participants (3,150 cognitively normal individuals (CN) and 4,266 AD patients), encompassing 10,764,329 SNPs. The male sex ratio was 56.3% for AD patients (mean age 70.1 years) and 60.7% for CN individuals (mean age 80.2 years).

Algorithm Implementation and Analysis

The Deep-Block framework employs a structured, three-stage process to analyze large-scale WGS data:

Stage 1: Segmentation of whole genome into LD blocks

Following QC, the WGS dataset was segmented into linkage disequilibrium (LD) blocks using Plink software. The parameters were set as follows: the LD measure was r^2 with a threshold of 0.9, window size of 50 variants, and maximum window physical size of 100 kilobases. LD blocks were then identified based on the genomic positions of SNPs and the extent of LD between adjacent SNPs. This configuration identified 30,218 LD Blocks, forming the basis for subsequent analyses.

Stage 2: Imputation of missing genotype data

Deep-Block utilizes machine learning approaches to impute missing genotype data within the LD blocks, a method supported by recent studies (An et al., 2023; Rubinacci et al., 2020; Shishegar et al., 2021). To

identify the most suitable imputer for imputing missing genotype data, preliminary experiments were conducted on the *APOE* gene region within the ADSP WGS dataset, assuming SNPs of this region are contained within LD blocks. The ADSP R3 WGS data, comprising 16,869 individuals, included 793 variants from the *APOE* gene region. After QC, the missing data proportion in this region was 4.14E-03. For the performance assessment of imputers, the missing rate was artificially increased to 8.70E-03. The modified dataset was then processed using the TopMed Imputer, establishing a benchmark for comparing the efficiency of other imputation methods. Several imputers were used: 1-NN, 5-NN, 10-NN, GAN, Iterative, MissForest (Stekhoven & Bühlmann, 2012), and Simple Imputer. All methods were applied to data with the same artificially increased proportion of missing genotype data to ensure consistent evaluation. The scikit-learn package (Pedregosa et al., 2011) was used for machine learning imputers and the GAIN package (Yoon et al., 2018) for the GAN Imputer.

The Simple Imputer utilizes mean, median, or mode imputation to fill missing values with the most representative statistic of the available data. The k-Nearest Neighbors (k-NN) Imputers (1-NN, 5-NN, and 10-NN) leverage data point similarity to impute missing values based on the nearest neighbors' values. The GAN Imputer uses Generative Adversarial Networks to produce synthetic data mimicking the original data distribution, thus imputing missing values. The Iterative Imputer employs a round-robin approach, modeling each feature with missing values as a function of other features stepwise, capturing complex interactions and dependencies. The MissForest Imputer utilizes a Random Forest approach, leveraging multiple decision trees to accurately predict missing values.

The performance of these methods was evaluated using five well-established metrics: accuracy, Root Mean Squared Error (RMSE), R-squared (R²), Mean Absolute Error (MAE), and Normalized RMSE (NRMSE). The accuracy quantifies the proportion of correctly imputed values, directly reflecting an imputer's performance. The RMSE measures the average magnitude of imputation errors, providing a straightforward accuracy metric. The R² indicates the proportion of variance in the original data explained by the imputed data, offering insights into the imputation method's ability to preserve data structure. The MAE calculates the average absolute error between imputed and actual values, presenting error distribution without directional bias. The NRMSE normalizes RMSE to the dataset range, facilitating the performance comparison across differently scaled datasets.

Stage 3: Identification of Key LD blocks and phenotype association

The final stage identifies key LD blocks as significantly associated with the AD phenotype using TabNet, a deep learning model optimized for efficient tabular data processing (Arik & Pfister, 2021). TabNet's architecture combines the interpretability of decision tree-based models with deep learning capabilities, featuring an encoder-decoder structure, feature transformers, and attentive transformers. TabNet's

encoder processes raw tabular data, selecting relevant features through a sequential multi-step procedure using feature transformers. These transformers apply non-linear transformations to enhance the model's learning capabilities. The attentive transformer, a key encoder component, employs the sparsemax normalization function to focus selectively on the most relevant features, optimizing model interpretability and efficiency. This stage uses TabNet to identify LD Blocks with high phenotypic relevance, focusing on features critically associated with AD. TabNet's decoder reconstructs features from the original dataset, identifying key features within the top LD Blocks. This process assigns Phenotype Influence Scores (PIS) to significant features, reflecting their phenotypic impact. The method integrates TabNet's feature importance metrics with the Mean Decrease Impurity (MDI) metric from Random Forest, offering a systematic approach to understanding genetic influences on phenotypic traits (Figure 1).

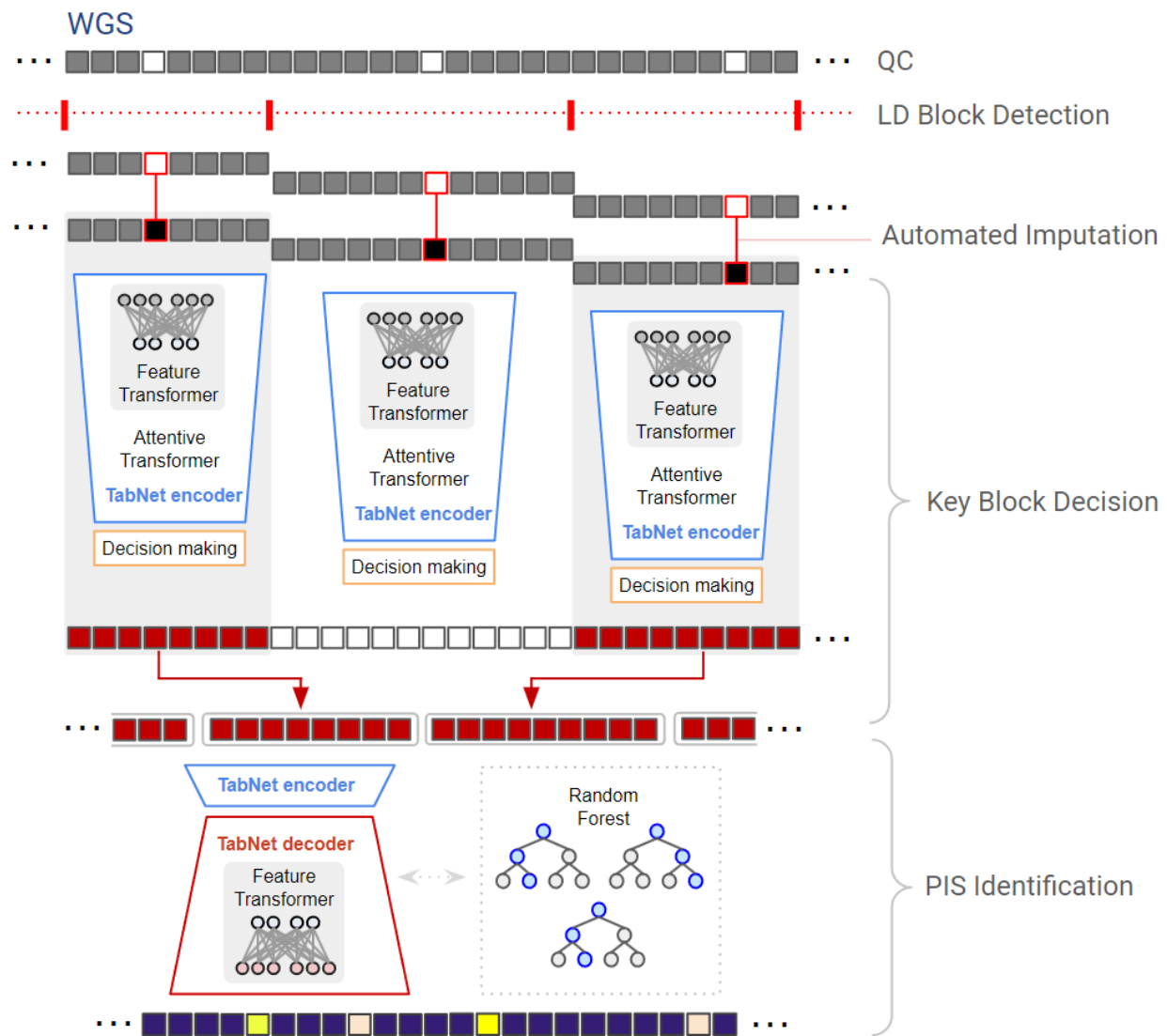


Figure 1. Overview of the Deep-Block Framework. This figure illustrates the sequential stages of the Deep-Block framework used in the analysis of large-scale whole genome sequencing (WGS) data for Alzheimer's disease (AD). The process initiates with the quality control procedure (QC) of WGS data, ensuring the integrity and reliability of the genetic information. Subsequently, the data is organized into Linkage Disequilibrium (LD) blocks, indicated by red dotted lines, which reflect the partitioning based on LD parameters. The next phase, Automated imputation, is visualized as various modules corresponding to different machine learning-based imputation techniques—each tasked with estimating and inputting missing genomic data. Following imputation, the TabNet encoder's role in decision-making is depicted, using feature transformers and attentive transformers to select and prioritize LD blocks that show significant associations with AD. The final element of the diagram focuses on the identification of the Phenotype Influence Score (PIS) using the TabNet decoder in conjunction with Random Forest metrics.

This approach combines the strengths of both metrics to identify the most significant AD-associated genetic markers, offering a robust method for detecting key genetic markers within LD blocks. The Phenotype Influence Score (PIS) is calculated using the following combined formula:

$$PIS_{Advanced_SWATj} = I \cdot M_{agg-b,j} + (1 - I) \cdot MDI_j$$

where I is an indicator variable that is automatically set to 1 when the TabNet model yields a higher predictive accuracy in phenotype-related classification using previously selected features, and is automatically set to 0 when the Random Forest algorithm shows superior performance in the same task. $M_{agg-b,j}$, the feature importance from the TabNet model, represents the aggregate feature importance mask for the j -th feature. The calculation uses the total number of decision steps (N), the learning rate at each decision step ($\eta_b[i]$), and a binary mask ($M_{b,j}[i]$) that is set to 1 if the j -th feature is utilized at the i -th decision step, and 0 otherwise. Here, D represents the total number of features. The corresponding formula is as follows:

$$M_{agg-b,j} = \frac{\sum_{i=1}^N \eta_b[i] M_{b,j}[i]}{\sum_{j=1}^D \sum_{i=1}^N \eta_b[i] M_{b,j}[i]}$$

MDI_j , the Mean Decrease Impurity from the Random Forest algorithm, quantifies the impurity reduction for a specific SNP (SNP_j). This calculation encompasses the total number of decision trees in the model (N_{trees}), each tree (t), and the node (i), using SNP_j for splitting, includes the number of samples at node i before the splitting (n_i^t) and the impurity reduction at this node ($\Delta i(s_i^t)$). The MDI is calculated as follows:

$$MDI_j = \frac{I}{N_{trees}} \sum_{t=1}^{N_{trees}} \sum_{i \in I_j^t} \frac{n_i^t}{n_{root}^t} \Delta i(s_i^t)$$

Results

This study analyzed large-scale WGS data from the ADSP, comprising 7,416 non-Hispanic White individuals (4,266 with Alzheimer's disease and 3,150 cognitively normal older adults). After quality control procedures, several imputation methods were comparatively evaluated: Simple, GAN, 1-NN, 5-NN, 10-NN, Iterative, MissForest, and TopMed Imputers. The assessment utilized metrics including accuracy, Root Mean Squared Error (RMSE), R-squared (R2), Mean Absolute Error (MAE), and Normalized RMSE (NRMSE).

The MissForest Imputer demonstrated superior performance among the machine learning-based methods, achieving the highest accuracy (0.999359), lowest RMSE (0.0039), and highest R2 (0.9993). The 5-NN and 10-NN Imputers also performed well, with accuracy rates of 0.999734 and 0.999626, R2 values of 0.9993, and RMSEs of 0.004 and 0.0041, respectively. The TopMed imputation server achieved an accuracy of 0.996416, RMSE of 0.0047, and R2 of 0.9081. While effective in reducing RMSE, it showed a lower capacity to capture dataset variance compared to the leading machine learning methods (Figure 2, Table 1).

Computation time for imputation methods was crucial due to the large-scale WGS data. Table 2 shows that the MissForest Imputer required up to 327 seconds for the largest block size, significantly longer than the 5-NN Imputer, which processed the same block in just over 50 milliseconds. Balancing imputation accuracy and processing efficiency, the 5-NN Imputer was selected as the most suitable imputation method for our dataset. This choice was based on its high accuracy and fast imputation capability.

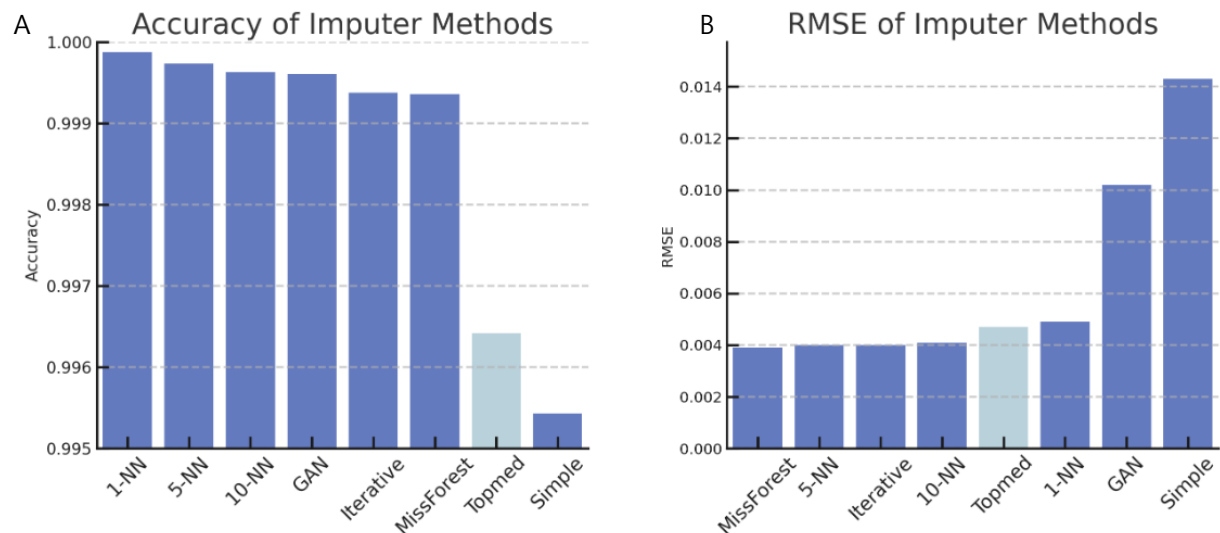


Figure 2. Comparative performance of imputation methods in WGS data. Fig. 2A illustrates the accuracy for several imputation methods, which reflects the proportion of correctly imputed genotypes to the total number of predictions made. A value closer to 1 denotes a higher rate of correct imputations. In this analysis, the 1-NN Imputer exhibits the highest accuracy, while the Simple Imputer shows the least accuracy, pointing to a greater discrepancy in its predictions. Fig. 2B displays the Root Mean Square Error (RMSE) across the imputation methods, a metric for quantifying the average errors in the predicted values. The lower the RMSE, the more accurate the imputation. Here, the MissForest Imputer emerges as the most accurate with the smallest RMSE, while the Simple Imputer displays the largest RMSE, indicative of lower accuracy. The results of the Topmed Imputer were not as pronounced, falling behind with lower accuracy and a higher RMSE than several other imputers.

IMPUTER METHOD	ACCURACY	RMSE	R2	MAE	NRMSE
1-NN IMPUTER	0.999875	0.0049	0.9989	0	0.0049
5-NN IMPUTER	0.999734	0.004	0.9993	0	0.004
10-NN IMPUTER	0.999626	0.0041	0.9993	0.0001	0.0041
GAN IMPUTER	0.999606	0.0102	0.9981	0.0002	0.0102
ITERATIVE IMPUTER	0.999373	0.004	0.9993	0.0001	0.004
MISSFOREST IMPUTER	0.999359	0.0039	0.9993	0	0.0039
TOPMED IMPUTER	0.996416	0.0047	0.9081	0.0002	0.0047
SIMPLE IMPUTER	0.995432	0.0143	0.9965	0.0006	0.0143

Table 1. Comparison of imputation efficacies of imputation methods. The table shows performance metrics for several imputation methods of missing genotypes. The accuracy measures the proportion of correctly imputed genotypes, where the 1-NN Imputer ranks the highest, suggesting the greatest precision in imputation among the methods. The Mean Squared Error (MSE) shows the MissForest Imputer as the most accurate, with the smallest values indicating minimal deviation from actual data. R-squared (R2) values for the 5-NN, 10-NN, Iterative, and MissForest Imputers indicate that these models account for a significant portion of the variance, suggesting a strong correlation with the observed data. The Mean Absolute Error (MAE) is lowest for the 1-NN, 5-NN, and MissForest Imputers, indicating higher precision. The Normalized Root Mean Squared Error (NRMSE) further confirms the MissForest Imputer's superior performance. Overall, the MissForest Imputer exhibits the highest precision in imputation of missing genotypes.

BLOCK SIZE	SIMPLE IMPUTER	5-NN IMPUTER	10-NN IMPUTER	MISSFOREST IMPUTER
40	0.0096	0.0076	0.0079	4.1889
80	0.0169	0.0122	0.0124	16.5368
120	0.0238	0.0174	0.0169	41.4959
160	0.0316	0.0236	0.0227	87.8926
200	0.0385	0.0257	0.0278	158.3003
300	0.0567	0.0387	0.0401	232.4933
400	0.0748	0.0506	0.0526	327.3879

Table 2. Comparison of computation time in imputing missing genotypes for imputation methods (measured in seconds). The table highlights the variation in processing time for imputation methods with increasing block sizes. While the Simple, 5-NN, and 10-NN Imputers show a gradual increase in processing time as block sizes increase, the MissForest Imputer exhibits a notably steep rise in computation time. This pronounced increase is especially significant for larger block sizes, where the MissForest Imputer's high accuracy is offset by its extensive processing time. Due to this substantial time consumption, the 5-NN Imputer, offering a balance between time efficiency and accuracy, was determined to be the most practical choice for our dataset.

LD blocks were determined using Plink, resulting in 30,218 LD blocks with an average size of 388 genetic variants. Genomic regions not covered by LD blocks, comprising only 0.19% of the genome, were excluded from the analysis due to their negligible size. Figure 3B demonstrates the correlation between the number of blocks per chromosome and chromosome length. TabNet was then applied to the LD blocks to assess phenotype prediction accuracy using a binary classification model. Blocks were ranked based on prediction accuracy, and as the number of blocks increased, the analysis showed that the number of important features with non-zero TabNet feature importance did not increase significantly. Figure 3A illustrates this by depicting genetic variants within the key blocks as blue bars. The analysis extended up to 1500 blocks, revealing a plateau in the count of important features, indicating that the critical variants for phenotype prediction were already captured within the initial top blocks. This suggests that further analysis beyond these blocks may not provide additional meaningful insights.

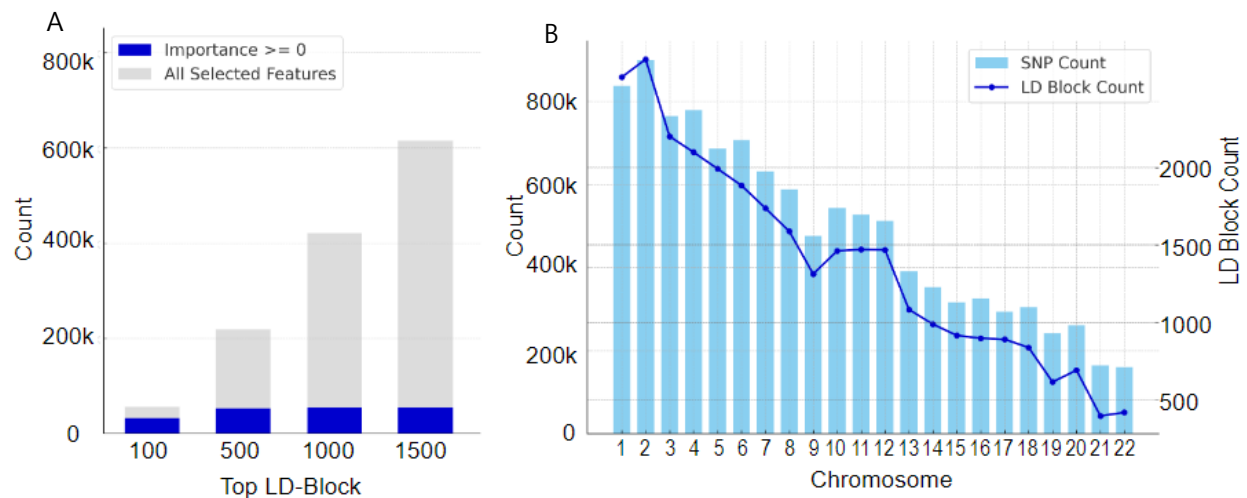


Figure 3. Feature importance and distribution across LD blocks in the ADSP WGS Dataset. In Fig. 3A, the analysis determined the feature importance using TabNet for the top 100 to 1500 LD blocks. The blue bars represent genetic variants within these blocks, where TabNet was assigned a feature importance greater than zero, indicating their relevance in phenotype prediction. The gray bars indicate all selected features, regardless of their importance score. The steady count of important features across increasing block ranks suggests that the most critical variants for phenotype prediction were concentrated in the top blocks. Fig 3B visualizes the distribution of LD blocks and SNP counts across chromosomes. The bar chart demonstrates that the number of LD blocks and SNPs is proportionate to the chromosome length, with larger chromosomes containing more blocks.

The study analyzed 54,949 genetic variants within the top 1500 blocks, calculating the PIS for each genetic variant using the methodology outlined in the Methods section. Table 3 presents the genetic variants with the highest PIS. The SNP rs429358 within the *APOE* gene on chromosome 19, a well-known AD risk SNP, demonstrated the highest importance score. Other high-importance SNPs, including rs11556505 in *TOMM40* and rs34342646 in *NECTIN2*, further emphasize the relevance of chromosome 19 in AD. The study also confirmed AD-associated SNPs, rs10414043, rs10119, and rs71352238, within the *APOC1* and *TOMM40* genes.

In addition to confirming the importance of known genes, this study identified novel high-importance SNPs not previously identified in genetic association studies for AD, particularly within the top 1500 LD Blocks examined. Notably, the study identified SNPs such as rs200986288 and rs199988716 (also known as rs62153752), along with genes like *LOC107984083* and *ANKRD30BL*, which have not been previously associated with AD in the literature. The study also identified previously unreported SNPs and their related genes in the context of AD, including rs66626994 (*APOC1P1*), rs73876031 (*FRG1*),

rs200608168 and rs201356360 (*ANKRD36*), and rs199988716 (*ANKRD36C*).

RANK	FEATURE	CHROMOSOME	POSITION	VARIATION	GENE	IMPORTANCE
1	rs429358	19	44908683	T/C	APOE	8.00E-04
2	rs11556505	19	44892886	C/T	TOMM40	7.23E-04
3	rs34342646	19	44884872	G/A	NECTIN2	6.85E-04
4	rs5117	19	44915532	T/C	APOC1	6.85E-04
5	rs483082	19	44912920	G/T	APOC1	6.09E-04
6	rs10414043	19	44912455	G/A	APOC1	5.43E-04
7	rs10119	19	44903415	G/A	TOMM40	5.11E-04
8	rs71352238	19	44891078	T/C	TOMM40	4.76E-04
9	rs6857	19	44888996	C/T	NECTIN2	4.63E-04
10	rs59007384	19	44893407	G/A G/T	TOMM40	4.57E-04
11	rs111789331	19	44923867	T/A		4.34E-04
12	rs12721046	19	44917996	G/A	APOC1	3.75E-04
13	rs4420638	19	44919688	A/G	APOC1	3.61E-04
14	rs283811	19	44885242	A/C A/G	NECTIN2	3.37E-04
15	.	N/A	N/A	N/A	N/A	3.06E-04
16	rs200986288	20	30185632	A/C A/T		2.87E-04
17	rs199988716	N/A	N/A	N/A	N/A	2.51E-04
18	rs202143966	9	63769614	T/C		2.49E-04
19	rs377656811	22	16427107	C/A C/T		2.45E-04
20	rs157582	19	44892961	C/T	TOMM40	2.37E-04
21	rs78790997	16	33736499	C/G	LOC107984083	2.28E-04
22	rs141490255	N/A	N/A	N/A	N/A	2.27E-04
23	rs75997270	4	189924739	C/A C/G	LOC728339	2.21E-04
24	rs75627662	19	44910318	C/T		2.20E-04
25	rs147747785	17	22046134	G/T		2.20E-04
26	rs1160985	19	44900154	C/T	TOMM40	2.18E-04
27	rs2075650	19	44892361	A/G	TOMM40	2.16E-04
28	rs202221379	17	22046133	G/T		2.14E-04
29	rs34404554	19	44892651	C/G	TOMM40	2.09E-04
30	rs769449	19	44906744	G/A	APOE	2.09E-04

Table 3. SNPs with highest phenotype influence scores (PIS) associated with AD. This table catalogs the top 30 single nucleotide polymorphisms (SNPs) ranked by PIS, derived from an extensive examination of 54,949 genetic variants within the top 1500 LD blocks using TabNet. The highest-scoring SNPs are predominantly located on chromosome 19, related to genes such as *APOE*, *APOC1*, *NECTIN2*, and *TOMM40*—well-established AD-associated genes. Additionally, this table includes novel findings, highlighting SNPs and genes previously unidentified in AD genetic association studies.

Figure 4 displays vertically aligned, symmetrical Manhattan plots (Miami plot) from two genetic association analysis methods: Deep-Block and Plink. The upper plot depicts the Phenotype Influence Scores (PIS) derived from Deep-Block, while the lower plot shows the statistical significance levels ($-\log_{10}$ P-value) obtained using Plink across all chromosomes. Both plots arrange chromosomes along the x-axis, offering a chromosomal position view of the analyzed genomic variants. Each plot highlights the top 40 genetic variants using color-coded dots: red for the top 1-10 genetic variants, blue for the top 11-20 genetic variants, green for the top 21-30 genetic variants, and gray for the top 31-40 genetic variants. Of note, the Deep-Block approach identified the same genetic variants that the Plink identified as well as novel genetic variants that the Plink could not identify.

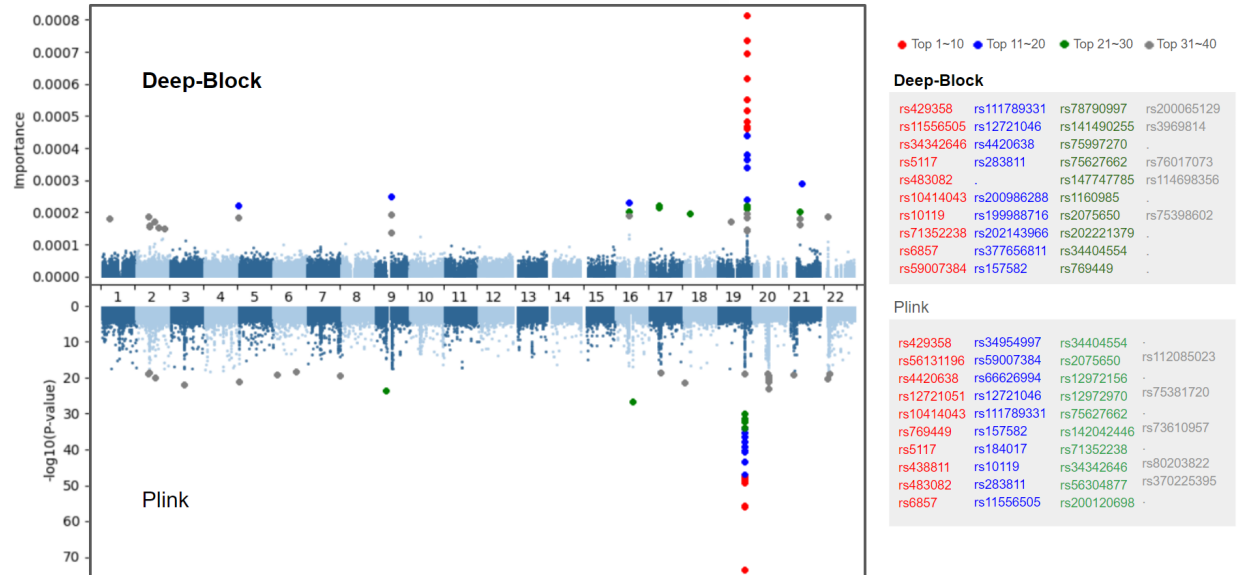


Figure 4. Miami plot showing comparison of association test statistics from between Deep-Block and Plink. The Miami plot displays Manhattan plots from two genetic association analysis methods, Deep-Block and Plink. The upper Manhattan plot shows genetic association test statistics from the Deep-Block framework, indicating the feature importance of single nucleotide polymorphisms (SNPs) throughout the genome, while the lower Manhattan plot presents the genetic association analysis results from Plink, reflecting the statistical significance of SNPs. The results correspond to the same genomic coordinates across the horizontal axis, which lists the chromosomes numerically. The significance of SNPs is visually encoded with colors: red dots mark the top 1-10 most important or significant SNPs, blue for the top 11-20, green for the top 21-30, and gray for the top 31-40. The side panel to the right of the plots enumerates SNPs within each color category.

Discussion

This study developed and applied the Deep-Block framework to large-scale WGS data from the ADSP to investigate the genetic basis of AD. The approach centered on utilizing LD blocks combined with automated imputation to improve the accuracy of genetic analysis. This methodological framework identified an array of genetic loci, including both previously identified known SNPs and newly identified novel SNPs.

The analysis revealed well-known associations between AD and several genetic loci (*APOE*, *TOMM40*, *NECTIN2*, and *APOC1*) and identified known variants (e.g., *APOE* rs429358 and rs769449). The findings align with previous genetic association studies, validating the effectiveness of the Deep-Block framework in identifying AD-associated genetic markers. In addition, the identification of both known and novel SNPs demonstrates this method's potential to broaden our understanding of AD's genetic factors.

To contextualize the Deep-Block findings, the identified loci were compared to established AD-linked genetic loci reported in the European Alzheimer's and Dementia Biobank (EADB) and the recently updated GWAS catalog (Lambert et al., 2023). The analysis aligned with 15 genetic loci from the GWAS catalog, distributed across various tiers (3 in Tier1, 1 in Tier2, 2 in Tier3, 2 unverified, and 3 in the 'Other' category), as detailed in Table S1. This comparison excluded the well-studied *APOE*, *TOMM40*/*APOC1*/*NECTIN2* loci to focus on other significant genetic associations. The results were also compared against databases referenced in the GWAS catalog, including IGAP2 (Kunkle et al., 2019), PGC1 (Jansen et al., 2019), IGAP2+UKB (Schwartzentruber et al., 2021), GR@ACE (de Rojas et al., 2021), PGC2 (Wightman et al., 2021), and EADB (Bellenguez et al., 2022). Table 4 highlights Tier1 genes *ABCA7*, *BIN1*, and *CR1*, identified in multiple studies. This confirms the Deep-Block method's relevance to known AD genetic markers and indicates both confirmatory and potentially novel genetic associations with the disease.

The analysis revealed previously unidentified genetic variants and genes (e.g., rs200986288, rs199988716, *LOC107984083*, and *ANKRD30BL*) associated with AD. These findings suggest additional genetic factors may influence AD, enhancing our understanding of the genetic architecture of AD. However, this study used sequencing data from non-Hispanic white individuals, which may limit the broad applicability of the findings. This limitation underscores the need for future studies involving more diverse participants to ensure broader relevance and applicability of findings across different populations.

Gene	Group	Study	GRCh38	rsID
ABCA7	Tier1	IGAP2	chr19:1050875:A:G	rs12151021

		PGC1	chr19:1039324:C:G	rs111278892
		IGAP2+UKB	chr19:1050875:A:G	rs12151021
		GR@ACE	chr19:1043639:C:T	rs3752231
		PGC2	chr19:1053525:C:G	rs3752241
		EADB	chr19:1050875:A:G	rs12151021
		Deep-Block	chr19:1045974:CCAGCC:CC	rs3764648
		BIN1	Tier1	IGAP2
PGC1	chr2:127133851:A:C			rs4663105
IGAP2+UKB	chr2:127135234:C:T			rs6733839
GR@ACE	chr2:127135234:C:T			rs6733839
PGC2	chr2:127133851:A:C			rs4663105
EADB	chr2:127135234:C:T			rs6733839
Deep-Block	chr2:127085030:A:G			rs10207708
CR1	Tier1	IGAP2	chr1:207577223:T:C	rs679515
		PGC1	chr1:207577223:T:C	rs679515
		IGAP2+UKB	chr1:207577223:T:C	rs679515
		GR@ACE	chr1:207629207:A:C	rs4844610
		PGC2	chr1:207577223:T:C	rs679515
		EADB	chr1:207518704:A:G	rs6656401
		Deep-Block	chr1:207594103:T:C	rs6697005

Table 4. Comparison between previously identified known genetic variants and findings from the Deep-Block framework. This table provides a cross-study comparison between known genetic variants and our findings associated with AD, focusing on the ABCA7, BIN1, and CR1 genes within the Tier 1 category. In particular, Deep-Block identified a SNP (rs150593211) in ABCA7, which is different from the rs12151021 variant commonly reported in previous research. Similarly, novel variants in the BIN1 and CR1 loci were also detected, showing the utility of Deep-Block in identifying genetic variants that may have been missed in traditional genetic association analyses.

Conclusion

This study developed and applied the Deep-Block AI framework to large-scale ADSP WGS data for genetic association analysis for AD. The approach involved segmenting the whole genome into LD blocks and applying automated imputation of missing genotypes for data preprocessing. The Deep-Block framework identified AD-associated genetic loci, including both previously identified and novel SNPs, leading to the identification of complex genetic patterns associated with AD that may have been overlooked in traditional genetic association methods and emphasizing the importance of advanced sequencing data analysis tools. Compared to traditional methods such as Plink and SKAT-O (Lee et al., 2012), Deep-Block uses LD structure and attention-based feature selection to analyze high-dimensional genomic data more comprehensively, potentially capturing genetic interactions that are not detected by conventional approaches. Unlike SWAT-CNN (Jo et al., 2022), which utilizes fixed genomic fragments, Deep-Block segments the genome based on LD-defined regions, thereby improving the identification

of biologically relevant patterns. This framework demonstrates its capability to analyze large-scale genomic data effectively and identify both known and novel genetic variants associated with Alzheimer's disease.

Funding

This research was partially supported by the Alzheimer's Association (AA) under the grant number AARG 22-974053 and the National Institutes of Health (NIH): P30 AG010133, P30 AG072976, R01 AG019771, R01 AG057739, U01 AG024904, R01 LM013463, R01 AG068193, T32 AG071444, U01 AG068057, U01AG072177, U19AG074879, R03AG063250, and R01 LM012535.

The Alzheimer's Disease Sequencing Project (ADSP) is comprised of two Alzheimer's Disease (AD) genetics consortia and three National Human Genome Research Institute (NHGRI) funded Large Scale Sequencing and Analysis Centers (LSAC). The two AD genetics consortia are the Alzheimer's Disease Genetics Consortium (ADGC) funded by NIA (U01 AG032984), and the Cohorts for Heart and Aging Research in Genomic Epidemiology (CHARGE) funded by NIA (R01 AG033193), the National Heart, Lung, and Blood Institute (NHLBI), other National Institute of Health (NIH) institutes and other foreign governmental and non-governmental organizations. The Discovery Phase analysis of sequence data is supported through UF1AG047133 (to Drs. Schellenberg, Farrer, Pericak-Vance, Mayeux, and Haines); U01AG049505 to Dr. Seshadri; U01AG049506 to Dr. Boerwinkle; U01AG049507 to Dr. Wijsman; and U01AG049508 to Dr. Goate and the Discovery Extension Phase analysis is supported through U01AG052411 to Dr. Goate, U01AG052410 to Dr. Pericak-Vance and U01 AG052409 to Drs. Seshadri and Fornage.

Sequencing for the Follow Up Study (FUS) is supported through U01AG057659 (to Drs. PericakVance, Mayeux, and Vardarajan) and U01AG062943 (to Drs. Pericak-Vance and Mayeux). Data generation and harmonization in the Follow-up Phase is supported by U54AG052427 (to Drs. Schellenberg and Wang). The FUS Phase analysis of sequence data is supported through U01AG058589 (to Drs. Destefano, Boerwinkle, De Jager, Fornage, Seshadri, and Wijsman), U01AG058654 (to Drs. Haines, Bush, Farrer, Martin, and Pericak-Vance), U01AG058635 (to Dr. Goate), RF1AG058066 (to Drs. Haines, Pericak-Vance, and Scott), RF1AG057519 (to Drs. Farrer and Jun), R01AG048927 (to Dr. Farrer), and RF1AG054074 (to Drs. Pericak-Vance and Beecham).

The ADGC cohorts include: Adult Changes in Thought (ACT) (U01 AG006781, U19 AG066567), the Alzheimer's Disease Research Centers (ADRC) (P30 AG062429, P30 AG066468, P30 AG062421, P30 AG066509, P30 AG066514, P30 AG066530, P30 AG066507, P30 AG066444, P30 AG066518, P30 AG066512, P30 AG066462, P30 AG072979, P30 AG072972, P30 AG072976, P30 AG072975, P30

AG072978, P30 AG072977, P30 AG066519, P30 AG062677, P30 AG079280, P30 AG062422, P30 AG066511, P30 AG072946, P30 AG062715, P30 AG072973, P30 AG066506, P30 AG066508, P30 AG066515, P30 AG072947, P30 AG072931, P30 AG066546, P20 AG068024, P20 AG068053, P20 AG068077, P20 AG068082, P30 AG072958, P30 AG072959), the Chicago Health and Aging Project (CHAP) (R01 AG11101, RC4 AG039085, K23 AG030944), Indiana Memory and Aging Study (IMAS) (R01 AG019771), Indianapolis Ibadan (R01 AG009956, P30 AG010133), the Memory and Aging Project (MAP) (R01 AG17917), Mayo Clinic (MAYO) (R01 AG032990, U01 AG046139, R01 NS080820, RF1 AG051504, P50 AG016574), Mayo Parkinson's Disease controls (NS039764, NS071674, 5RC2HG005605), University of Miami (R01 AG027944, R01 AG028786, R01 AG019085, IIRG09133827, A2011048), the Multi-Institutional Research in Alzheimer's Genetic Epidemiology Study (MIRAGE) (R01 AG09029, R01 AG025259), the National Centralized Repository for Alzheimer's Disease and Related Dementias (NCRAD) (U24 AG021886), the National Institute on Aging Late Onset Alzheimer's Disease Family Study (NIA-LOAD) (U24 AG056270), the Religious Orders Study (ROS) (P30 AG10161, R01 AG15819), the Texas Alzheimer's Research and Care Consortium (TARCC) (funded by the Darrell K Royal Texas Alzheimer's Initiative), Vanderbilt University/Case Western Reserve University (VAN/CWRU) (R01 AG019757, R01 AG021547, R01 AG027944, R01 AG028786, P01 NS026630, and Alzheimer's Association), the Washington Heights-Inwood Columbia Aging Project (WHICAP) (RF1 AG054023), the University of Washington Families (VA Research Merit Grant, NIA: P50AG005136, R01AG041797, NINDS: R01NS069719), the Columbia University Hispanic Estudio Familiar de Influencia Genetica de Alzheimer (EFIGA) (RF1 AG015473), the University of Toronto (UT) (funded by Wellcome Trust, Medical Research Council, Canadian Institutes of Health Research), and Genetic Differences (GD) (R01 AG007584). The CHARGE cohorts are supported in part by National Heart, Lung, and Blood Institute (NHLBI) infrastructure grant HL105756 (Psaty), RC2HL102419 (Boerwinkle) and the neurology working group is supported by the National Institute on Aging (NIA) R01 grant AG033193.

The CHARGE cohorts participating in the ADSP include the following: Austrian Stroke Prevention Study (ASPS), ASPS-Family study, and the Prospective Dementia Registry-Austria (ASPS/PRODEM-Aus), the Atherosclerosis Risk in Communities (ARIC) Study, the Cardiovascular Health Study (CHS), the Erasmus Rucphen Family Study (ERF), the Framingham Heart Study (FHS), and the Rotterdam Study (RS). ASPS is funded by the Austrian Science Fond (FWF) grant number P20545-P05 and P13180 and the Medical University of Graz. The ASPS-Fam is funded by the Austrian Science Fund (FWF) project I904), the EU Joint Programme – Neurodegenerative Disease Research (JPND) in frame of the BRIDGET project (Austria, Ministry of Science) and the Medical University of Graz and the Steiermärkische Krankenanstalten Gesellschaft. PRODEM-Austria is supported by the Austrian Research Promotion agency (FFG) (Project No. 827462) and by the Austrian National Bank (Anniversary Fund, project 15435. ARIC research is

carried out as a collaborative study supported by NHLBI contracts (HHSN268201100005C, HHSN268201100006C, HHSN268201100007C, HHSN268201100008C, HHSN268201100009C, HHSN268201100010C, HHSN268201100011C, and HHSN268201100012C). Neurocognitive data in ARIC is collected by U01 2U01HL096812, 2U01HL096814, 2U01HL096899, 2U01HL096902, 2U01HL096917 from the NIH (NHLBI, NINDS, NIA and NIDCD), and with previous brain MRI examinations funded by R01-HL70825 from the NHLBI. CHS research was supported by contracts HHSN268201200036C, HHSN268200800007C, N01HC55222, N01HC85079, N01HC85080, N01HC85081, N01HC85082, N01HC85083, N01HC85086, and grants U01HL080295 and U01HL130114 from the NHLBI with additional contribution from the National Institute of Neurological Disorders and Stroke (NINDS). Additional support was provided by R01AG023629, R01AG15928, and R01AG20098 from the NIA. FHS research is supported by NHLBI contracts N01-HC-25195 and HHSN268201500001I. This study was also supported by additional grants from the NIA (R01s AG054076, AG049607 and AG033040 and NINDS (R01 NS017950). The ERF study as a part of EUROSPAN (European Special Populations Research Network) was supported by European Commission FP6 STRP grant number 018947 (LSHG-CT-2006-01947) and also received funding from the European Community's Seventh Framework Programme (FP7/2007-2013)/grant agreement HEALTH-F4- 2007-201413 by the European Commission under the programme "Quality of Life and Management of the Living Resources" of 5th Framework Programme (no. QL62-CT-2002- 01254). High-throughput analysis of the ERF data was supported by a joint grant from the Netherlands Organization for Scientific Research and the Russian Foundation for Basic Research (NWO-RFBR 047.017.043). The Rotterdam Study is funded by Erasmus Medical Center and Erasmus University, Rotterdam, the Netherlands Organization for Health Research and Development (ZonMw), the Research Institute for Diseases in the Elderly (RIDE), the Ministry of Education, Culture and Science, the Ministry for Health, Welfare and Sports, the European Commission (DG XII), and the municipality of Rotterdam. Genetic data sets are also supported by the Netherlands Organization of Scientific Research NWO Investments (175.010.2005.011, 911-03-012), the Genetic Laboratory of the Department of Internal Medicine, Erasmus MC, the Research Institute for Diseases in the Elderly (014-93-015; RIDE2), and the Netherlands Genomics Initiative (NGI)/Netherlands Organization for Scientific Research (NWO) Netherlands Consortium for Healthy Aging (NCHA), project 050-060-810. All studies are grateful to their participants, faculty and staff. The content of these manuscripts is solely the responsibility of the authors and does not necessarily represent the official views of the National Institutes of Health or the U.S. Department of Health and Human Services.

The FUS cohorts include: the Alzheimer's Disease Research Centers (ADRC) (P30 AG062429, P30 AG066468, P30 AG062421, P30 AG066509, P30 AG066514, P30 AG066530, P30 AG066507, P30 AG066444, P30 AG066518, P30 AG066512, P30 AG066462, P30 AG072979, P30 AG072972, P30

AG072976, P30 AG072975, P30 AG072978, P30 AG072977, P30 AG066519, P30 AG062677, P30 AG079280, P30 AG062422, P30 AG066511, P30 AG072946, P30 AG062715, P30 AG072973, P30 AG066506, P30 AG066508, P30 AG066515, P30 AG072947, P30 AG072931, P30 AG066546, P20 AG068024, P20 AG068053, P20 AG068077, P20 AG068082, P30 AG072958, P30 AG072959), Alzheimer's Disease Neuroimaging Initiative (ADNI) (U19AG024904), Amish Protective Variant Study (RF1AG058066), Cache County Study (R01AG11380, R01AG031272, R01AG21136, RF1AG054052), Case Western Reserve University Brain Bank (CWRUBB) (P50AG008012), Case Western Reserve University Rapid Decline (CWRURD) (RF1AG058267, NU38CK000480), CubanAmerican Alzheimer's Disease Initiative (CuADI) (3U01AG052410), Estudio Familiar de Influencia Genetica en Alzheimer (EFIGA) (5R37AG015473, RF1AG015473, R56AG051876), Genetic and Environmental Risk Factors for Alzheimer Disease Among African Americans Study (GenerAAtions) (2R01AG09029, R01AG025259, 2R01AG048927), Gwangju Alzheimer and Related Dementias Study (GARD) (U01AG062602), Hillblom Aging Network (2014-A-004-NET, R01AG032289, R01AG048234), Hussman Institute for Human Genomics Brain Bank (HIHGBB) (R01AG027944, Alzheimer's Association "Identification of Rare Variants in Alzheimer Disease"), Ibadan Study of Aging (IBADAN) (5R01AG009956), Longevity Genes Project (LGP) and LonGenity (R01AG042188, R01AG044829, R01AG046949, R01AG057909, R01AG061155, P30AG038072), Mexican Health and Aging Study (MHAS) (R01AG018016), Multi-Institutional Research in Alzheimer's Genetic Epidemiology (MIRAGE) (2R01AG09029, R01AG025259, 2R01AG048927), Northern Manhattan Study (NOMAS) (R01NS29993), Peru Alzheimer's Disease Initiative (PeADI) (RF1AG054074), Puerto Rican 1066 (PR1066) (Wellcome Trust (GR066133/GR080002), European Research Council (340755)), Puerto Rican Alzheimer Disease Initiative (PRADI) (RF1AG054074), Reasons for Geographic and Racial Differences in Stroke (REGARDS) (U01NS041588), Research in African American Alzheimer Disease Initiative (REAAADI) (U01AG052410), the Religious Orders Study (ROS) (P30 AG10161, P30 AG72975, R01 AG15819, R01 AG42210), the RUSH Memory and Aging Project (MAP) (R01 AG017917, R01 AG42210Stanford Extreme Phenotypes in AD (R01AG060747), University of Miami Brain Endowment Bank (MBB), University of Miami/Case Western/North Carolina A&T African American (UM/CASE/NCAT) (U01AG052410, R01AG028786), and Wisconsin Registry for Alzheimer's Prevention (WRAP) (R01AG027161 and R01AG054047).

The four LSACs are: the Human Genome Sequencing Center at the Baylor College of Medicine (U54 HG003273), the Broad Institute Genome Center (U54HG003067), The American Genome Center at the Uniformed Services University of the Health Sciences (U01AG057659), and the Washington University Genome Institute (U54HG003079). Genotyping and sequencing for the ADSP FUS is also conducted at John P. Hussman Institute for Human Genomics (HIHG) Center for Genome Technology (CGT).

Biological samples and associated phenotypic data used in primary data analyses were stored at Study Investigators institutions, and at the National Centralized Repository for Alzheimer's Disease and Related Dementias (NCRAD, U24AG021886) at Indiana University funded by NIA. Associated Phenotypic Data used in primary and secondary data analyses were provided by Study Investigators, the NIA funded Alzheimer's Disease Centers (ADCs), and the National Alzheimer's Coordinating Center (NACC, U24AG072122) and the National Institute on Aging Genetics of Alzheimer's Disease Data Storage Site (NIAGADS, U24AG041689) at the University of Pennsylvania, funded by NIA. Harmonized phenotypes were provided by the ADSP Phenotype Harmonization Consortium (ADSP-PHC), funded by NIA (U24 AG074855, U01 AG068057 and R01 AG059716) and Ultrascale Machine Learning to Empower Discovery in Alzheimer's Disease Biobanks (AI4AD, U01 AG068057). This research was supported in part by the Intramural Research Program of the National Institutes of Health, National Library of Medicine. Contributors to the Genetic Analysis Data included Study Investigators on projects that were individually funded by NIA, and other NIH institutes, and by private U.S. organizations, or foreign governmental or nongovernmental organizations.

The ADSP Phenotype Harmonization Consortium (ADSP-PHC) is funded by NIA (U24 AG074855, U01 AG068057 and R01 AG059716). The harmonized cohorts within the ADSP-PHC include: the Anti-Amyloid Treatment in Asymptomatic Alzheimer's study (A4 Study), a secondary prevention trial in preclinical Alzheimer's disease, aiming to slow cognitive decline associated with brain amyloid accumulation in clinically normal older individuals. The A4 Study is funded by a public-private-philanthropic partnership, including funding from the National Institutes of Health-National Institute on Aging, Eli Lilly and Company, Alzheimer's Association, Accelerating Medicines Partnership, GHR Foundation, an anonymous foundation and additional private donors, with in-kind support from Avid and Cogstate. The companion observational Longitudinal Evaluation of Amyloid Risk and Neurodegeneration (LEARN) Study is funded by the Alzheimer's Association and GHR Foundation. The A4 and LEARN Studies are led by Dr. Reisa Sperling at Brigham and Women's Hospital, Harvard Medical School and Dr. Paul Aisen at the Alzheimer's Therapeutic Research Institute (ATRI), University of Southern California. The A4 and LEARN Studies are coordinated by ATRI at the University of Southern California, and the data are made available through the Laboratory for Neuro Imaging at the University of Southern California. The participants screening for the A4 Study provided permission to share their de-identified data in order to advance the quest to find a successful treatment for Alzheimer's disease. We would like to acknowledge the dedication of all the participants, the site personnel, and all of the partnership team members who continue to make the A4 and LEARN Studies possible. The complete A4 Study Team list is available on: a4study.org/a4-study-team; the Adult Changes in Thought study (ACT), U01 AG006781, U19 AG066567; Alzheimer's Disease Neuroimaging Initiative (ADNI): Data collection and sharing for this project was

funded by the Alzheimer's Disease Neuroimaging Initiative (ADNI) (National Institutes of Health Grant U01 AG024904) and DOD ADNI (Department of Defense award number W81XWH-12-2-0012). ADNI is funded by the National Institute on Aging, the National Institute of Biomedical Imaging and Bioengineering, and through generous contributions from the following: AbbVie, Alzheimer's Association; Alzheimer's Drug Discovery Foundation; Araclon Biotech; BioClinica, Inc.; Biogen; Bristol-Myers Squibb Company; CereSpir, Inc.; Cogstate; Eisai Inc.; Elan Pharmaceuticals, Inc.; Eli Lilly and Company; EuroImmun; F. Hoffmann-La Roche Ltd and its affiliated company Genentech, Inc.; Fujirebio; GE Healthcare; IXICO Ltd.; Janssen Alzheimer Immunotherapy Research & Development, LLC.; Johnson & Johnson Pharmaceutical Research & Development LLC.; Lumosity; Lundbeck; Merck & Co., Inc.; Meso Scale Diagnostics, LLC.; NeuroRx Research; Neurotrack Technologies; Novartis Pharmaceuticals Corporation; Pfizer Inc.; Piramal Imaging; Servier; Takeda Pharmaceutical Company; and Transition Therapeutics. The Canadian Institutes of Health Research is providing funds to support ADNI clinical sites in Canada. Private sector contributions are facilitated by the Foundation for the National Institutes of Health (www.fnih.org). The grantee organization is the Northern California Institute for Research and Education, and the study is coordinated by the Alzheimer's Therapeutic Research Institute at the University of Southern California. ADNI data are disseminated by the Laboratory for Neuro Imaging at the University of Southern California; Estudio Familiar de Influencia Genética en Alzheimer (EFIGA): 5R37AG015473, RF1AG015473, R56AG051876; Memory & Aging Project at Knight Alzheimer's Disease Research Center (MAP at Knight ADRC): The Memory and Aging Project at the Knight-ADRC (Knight-ADRC). This work was supported by the National Institutes of Health (NIH) grants R01AG064614, R01AG044546, RF1AG053303, RF1AG058501, U01AG058922 and R01AG064877 to Carlos Cruchaga. The recruitment and clinical characterization of research participants at Washington University was supported by NIH grants P30AG066444, P01AG03991, and P01AG026276. Data collection and sharing for this project was supported by NIH grants RF1AG054080, P30AG066462, R01AG064614 and U01AG052410. We thank the contributors who collected samples used in this study, as well as patients and their families, whose help and participation made this work possible. This work was supported by access to equipment made possible by the Hope Center for Neurological Disorders, the Neurogenomics and Informatics Center (NGI: <https://neurogenomics.wustl.edu/>) and the Departments of Neurology and Psychiatry at Washington University School of Medicine; National Alzheimer's Coordinating Center (NACC): The NACC database is funded by NIA/NIH Grant U24 AG072122. NACC data are contributed by the NIA-funded ADRCs: P30 AG062429 (PI James Brewer, MD, PhD), P30 AG066468 (PI Oscar Lopez, MD), P30 AG062421 (PI Bradley Hyman, MD, PhD), P30 AG066509 (PI Thomas Grabowski, MD), P30 AG066514 (PI Mary Sano, PhD), P30 AG066530 (PI Helena Chui, MD), P30 AG066507 (PI Marilyn Albert, PhD), P30 AG066444 (PI John Morris, MD), P30 AG066518 (PI Jeffrey Kaye, MD), P30 AG066512 (PI

Thomas Wisniewski, MD), P30 AG066462 (PI Scott Small, MD), P30 AG072979 (PI David Wolk, MD), P30 AG072972 (PI Charles DeCarli, MD), P30 AG072976 (PI Andrew Saykin, PsyD), P30 AG072975 (PI David Bennett, MD), P30 AG072978 (PI Neil Kowall, MD), P30 AG072977 (PI Robert Vassar, PhD), P30 AG066519 (PI Frank LaFerla, PhD), P30 AG062677 (PI Ronald Petersen, MD, PhD), P30 AG079280 (PI Eric Reiman, MD), P30 AG062422 (PI Gil Rabinovici, MD), P30 AG066511 (PI Allan Levey, MD, PhD), P30 AG072946 (PI Linda Van Eldik, PhD), P30 AG062715 (PI Sanjay Asthana, MD, FRCP), P30 AG072973 (PI Russell Swerdlow, MD), P30 AG066506 (PI Todd Golde, MD, PhD), P30 AG066508 (PI Stephen Strittmatter, MD, PhD), P30 AG066515 (PI Victor Henderson, MD, MS), P30 AG072947 (PI Suzanne Craft, PhD), P30 AG072931 (PI Henry Paulson, MD, PhD), P30 AG066546 (PI Sudha Seshadri, MD), P20 AG068024 (PI Erik Roberson, MD, PhD), P20 AG068053 (PI Justin Miller, PhD), P20 AG068077 (PI Gary Rosenberg, MD), P20 AG068082 (PI Angela Jefferson, PhD), P30 AG072958 (PI Heather Whitson, MD), P30 AG072959 (PI James Leverenz, MD); National Institute on Aging Alzheimer's Disease Family Based Study (NIA-AD FBS): U24 AG056270; Religious Orders Study (ROS): P30AG10161, R01AG15819, R01AG42210; Memory and Aging Project (MAP - Rush): R01AG017917, R01AG42210; Minority Aging Research Study (MARS): R01AG22018, R01AG42210; Washington Heights/Inwood Columbia Aging Project (WHICAP): RF1 AG054023; and Wisconsin Registry for Alzheimer's Prevention (WRAP): R01AG027161 and R01AG054047. Additional acknowledgments include the National Institute on Aging Genetics of Alzheimer's Disease Data Storage Site (NIAGADS, U24AG041689) at the University of Pennsylvania, funded by NIA.

Conflict of Interest: Dr. Saykin receives support from multiple NIH grants (P30 AG010133, P30 AG072976, R01 AG019771, R01 AG057739, U19 AG024904, R01 LM013463, R01 AG068193, T32 AG071444, U01 AG068057, U01 AG072177, and U19 AG074879). He has also received support from Avid Radiopharmaceuticals, a subsidiary of Eli Lilly (in kind contribution of PET tracer precursor) and participated in Scientific Advisory Boards (Bayer Oncology, Eisai, Novo Nordisk, and Siemens Medical Solutions USA, Inc) and an Observational Study Monitoring Board (MESA, NIH NHLBI), as well as External Advisory Committees for multiple NIA grants. He also serves as Editor-in-Chief of Brain Imaging and Behavior, a Springer-Nature Journal. The other authors declare no conflict of interest.

References

- An, U., Pazokitoroudi, A., Alvarez, M., Huang, L., Bacanu, S., Schork, A. J., Kendler, K., Pajukanta, P., Flint, J., Zaitlen, N., Cai, N., Dahl, A., & Sankararaman, S. (2023). Deep learning-based phenotype imputation on population-scale biobank data increases genetic discoveries. *Nature Genetics*, 55(12), 2269-2276. <https://doi.org/10.1038/s41588-023-01558-w>
- Arik, S. Ö., & Pfister, T. (2021). Tabnet: Attentive interpretable tabular learning. Proceedings of the AAAI Conference on Artificial Intelligence,

- Avsec, Ž., Agarwal, V., Visentin, D., Ledsam, J. R., Grabska-Barwinska, A., Taylor, K. R., Assael, Y., Jumper, J., Kohli, P., & Kelley, D. R. (2021). Effective gene expression prediction from sequence by integrating long-range interactions. *Nature Methods*, *18*(10), 1196-1203. <https://doi.org/10.1038/s41592-021-01252-x>
- Bellenguez, C., Küçükali, F., Jansen, I. E., Kleineidam, L., Moreno-Grau, S., Amin, N., Naj, A. C., Campos-Martin, R., Grenier-Boley, B., Andrade, V., Holmans, P. A., Boland, A., Damotte, V., van der Lee, S. J., Costa, M. R., Kuulasmaa, T., Yang, Q., de Rojas, I., Bis, J. C., . . . Sánchez-Arjona, M. B. (2022). New insights into the genetic etiology of Alzheimer's disease and related dementias. *Nature Genetics*, *54*(4), 412-436. <https://doi.org/10.1038/s41588-022-01024-z>
- Berson, E., Sreenivas, A., Phongpreecha, T., Perna, A., Grandi, F. C., Xue, L., Ravindra, N. G., Payrovnaziri, N., Mataraso, S., Kim, Y., Espinosa, C., Chang, A. L., Becker, M., Montine, K. S., Fox, E. J., Chang, H. Y., Corces, M. R., Aghaeepour, N., & Montine, T. J. (2023). Whole genome deconvolution unveils Alzheimer's resilient epigenetic signature. *Nature communications*, *14*(1), 4947. <https://doi.org/10.1038/s41467-023-40611-4>
- Bettencourt, C., Skene, N., Bandres-Ciga, S., Anderson, E., Winchester, L. M., Foote, I. F., Schwartzentruber, J., Botia, J. A., Nalls, M., Singleton, A., Schilder, B. M., Humphrey, J., Marzi, S. J., Toomey, C. E., Kleiflat, A. A., Harshfield, E. L., Garfield, V., Sandor, C., Keat, S., . . . Llewellyn, D. J. Artificial intelligence for dementia genetics and omics. *Alzheimer's & dementia*, *n/a*(n/a). <https://doi.org/https://doi.org/10.1002/alz.13427>
- de Rojas, I., Moreno-Grau, S., Tesi, N., Grenier-Boley, B., Andrade, V., Jansen, I. E., Pedersen, N. L., Stringa, N., Zettergren, A., Hernández, I., Montreal, L., Antúnez, C., Antonell, A., Tankard, R. M., Bis, J. C., Sims, R., Bellenguez, C., Quintela, I., González-Perez, A., . . . contributors, E. (2021). Common variants in Alzheimer's disease and risk stratification by polygenic risk scores. *Nature Communications*, *12*(1), 3417. <https://doi.org/10.1038/s41467-021-22491-8>
- Eraslan, G., Avsec, Ž., Gagneur, J., & Theis, F. J. (2019). Deep learning: new computational modelling techniques for genomics. *Nature Reviews Genetics*, *20*(7), 389-403. <https://doi.org/10.1038/s41576-019-0122-6>
- Fujiwara, T., Kwon, O.-H., & Ma, K. L. (2020). Supporting Analysis of Dimensionality Reduction Results With Contrastive Learning. *Ieee Transactions on Visualization and Computer Graphics*. <https://doi.org/10.1109/tvcg.2019.2934251>
- Jansen, I. E., Savage, J. E., Watanabe, K., Bryois, J., Williams, D. M., Steinberg, S., Sealock, J., Karlsson, I. K., Hägg, S., Athanasiu, L., Voyle, N., Proitsi, P., Witoelar, A., Stringer, S., Aarsland, D., Almdahl, I. S., Andersen, F., Bergh, S., Bettella, F., . . . Posthuma, D. (2019). Genome-wide meta-analysis identifies new loci and functional pathways influencing Alzheimer's disease risk. *Nature Genetics*, *51*(3), 404-413. <https://doi.org/10.1038/s41588-018-0311-9>
- Jo, T., Kim, J., Bice, P., Huynh, K., Wang, T., Arnold, M., Meikle, P. J., Giles, C., Kaddurah-Daouk, R., Saykin, A. J., & Nho, K. (2023). Circular-SWAT for deep learning based diagnostic classification of

- Alzheimer's disease: application to metabolome data. *EBioMedicine*, *97*, 104820. <https://doi.org/10.1016/j.ebiom.2023.104820>
- Jo, T., Nho, K., Bice, P., Saykin, A. J., & Initiative, A. S. D. N. (2022). Deep learning-based identification of genetic variants: application to Alzheimer's disease classification. *Briefings in Bioinformatics*, *23*(2), bbac022.
- Karczewski, K. J., & Snyder, M. P. (2018). Integrative omics for health and disease. *Nature Reviews Genetics*, *19*(5), 299-310. <https://doi.org/10.1038/nrg.2018.4>
- Konietschke, F., Schwab, K., & Pauly, M. (2021). Small sample sizes: A big data problem in high-dimensional data analysis. *Stat Methods Med Res*, *30*(3), 687-701. <https://doi.org/10.1177/0962280220970228>
- Kunkle, B. W., Grenier-Boley, B., Sims, R., Bis, J. C., Damotte, V., Naj, A. C., Boland, A., Vronskaya, M., van der Lee, S. J., Amlie-Wolf, A., Bellenguez, C., Frizatti, A., Chouraki, V., Martin, E. R., Sleegers, K., Badarinarayan, N., Jakobsdottir, J., Hamilton-Nelson, K. L., Moreno-Grau, S., . . . Lieberman, A. P. (2019). Genetic meta-analysis of diagnosed Alzheimer's disease identifies new risk loci and implicates A β , tau, immunity and lipid processing. *Nature Genetics*, *51*(3), 414-430. <https://doi.org/10.1038/s41588-019-0358-2>
- Lambert, J.-C., Ramirez, A., Grenier-Boley, B., & Bellenguez, C. (2023). Step by step: towards a better understanding of the genetic architecture of Alzheimer's disease. *Molecular Psychiatry*, *28*(7), 2716-2727. <https://doi.org/10.1038/s41380-023-02076-1>
- Lee, S., Wu, M. C., & Lin, X. (2012). Optimal tests for rare variant effects in sequencing association studies. *Biostatistics*, *13*(4), 762-775. <https://doi.org/10.1093/biostatistics/kxs014>
- Novakovsky, G., Dexter, N., Libbrecht, M. W., Wasserman, W. W., & Mostafavi, S. (2023). Obtaining genetics insights from deep learning via explainable artificial intelligence. *Nature Reviews Genetics*, *24*(2), 125-137. <https://doi.org/10.1038/s41576-022-00532-2>
- Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M., Prettenhofer, P., Weiss, R., & Dubourg, V. (2011). Scikit-learn: Machine learning in Python. *the Journal of machine Learning research*, *12*, 2825-2830.
- Rubinacci, S., Delaneau, O., & Marchini, J. (2020). Genotype Imputation Using the Positional Burrows Wheeler Transform. *Plos Genetics*. <https://doi.org/10.1371/journal.pgen.1009049>
- Schwartzentruber, J., Cooper, S., Liu, J. Z., Barrio-Hernandez, I., Bello, E., Kumasaka, N., Young, A. M., Franklin, R. J., Johnson, T., & Estrada, K. (2021). Genome-wide meta-analysis, fine-mapping and integrative prioritization implicate new Alzheimer's disease risk genes. *Nature Genetics*, *53*(3), 392-402.
- Shetta, O., & Niranjana, M. (2020). Robust Subspace Methods for Outlier Detection in Genomic Data Circumvents the Curse of Dimensionality. *Royal Society Open Science*. <https://doi.org/10.1098/rsos.190714>
- Shigemizu, D., Akiyama, S., Sukanuma, M., Furutani, M., Yamakawa, A., Nakano, Y., Ozaki, K., & Niida, S.

- (2023). Classification and deep-learning–based prediction of Alzheimer disease subtypes by using genomic data. *Translational psychiatry*, 13(1), 232. <https://doi.org/10.1038/s41398-023-02531-1>
- Shishegar, R., Cox, T., Rolls, D. A., Bourgeat, P., Dore, V., Lamb, F., Robertson, J., Laws, S. M., Porter, T., Fripp, J., Tosun, D., Maruff, P., Savage, G., Rowe, C. C., Masters, C. L., Weiner, M. W., Villemagne, V. L., & Burnham, S. (2021). Using Imputation to Provide Harmonized Longitudinal Measures of Cognition Across AIBL and ADNI. *Scientific Reports*. <https://doi.org/10.1038/s41598-021-02827-6>
- Stekhoven, D. J., & Bühlmann, P. (2012). MissForest—non-parametric missing value imputation for mixed-type data. *Bioinformatics*, 28(1), 112-118.
- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, Ł., & Polosukhin, I. (2017). Attention is all you need. *Advances in neural information processing systems*, 30.
- Vogelstein, J. T., Bridgeford, E., Tang, M., Zheng, D., Douville, C., Burns, R., & Maggioni, M. (2021). Supervised Dimensionality Reduction for Big Data. *Nature communications*. <https://doi.org/10.1038/s41467-021-23102-2>
- Wightman, D. P., Jansen, I. E., Savage, J. E., Shadrin, A. A., Bahrami, S., Holland, D., Rongve, A., Børte, S., Winsvold, B. S., Drange, O. K., Martinsen, A. E., Skogholt, A. H., Willer, C., Bråthen, G., Bosnes, I., Nielsen, J. B., Fritsche, L. G., Thomas, L. F., Pedersen, L. M., . . . andMe Research, T. (2021). A genome-wide association study with 1,126,563 individuals identifies new risk loci for Alzheimer’s disease. *Nature Genetics*, 53(9), 1276-1282. <https://doi.org/10.1038/s41588-021-00921-z>
- Yoon, J., Jordon, J., & Schaar, M. (2018). Gain: Missing data imputation using generative adversarial nets. International conference on machine learning,
- Zhou, J., & Troyanskaya, O. G. (2015). Predicting effects of noncoding variants with deep learning–based sequence model. *Nature Methods*, 12(10), 931-934. <https://doi.org/10.1038/nmeth.3547>