

1           Leveraging functional annotations to map rare  
2 variants associated with Alzheimer's disease with  
3 **gruyere**

4           Anjali Das<sup>1,2</sup>, Chirag Lakhani<sup>2</sup>, Chloé Terwagne<sup>3</sup>,  
5 Jui-Shan T. Lin<sup>2</sup>, Tatsuhiko Naito<sup>2,4</sup>, Towfique Raj<sup>4</sup>,  
6 David A. Knowles<sup>1,2,5,6\*</sup>

7 <sup>1</sup>Computer Science, Columbia University, New York, NY, USA.

8 <sup>2</sup>New York Genome Center, New York, NY, USA.

9 <sup>3</sup>Francis Crick Institute, London, United Kingdom.

10 <sup>4</sup>Neuroscience, Icahn School of Medicine, Mount Sinai, New York, NY,  
11 USA.

12 <sup>5</sup>Systems Biology, Columbia University, New York, NY, USA.

13 <sup>6</sup>Data Science Institute, Columbia University, New York, NY, USA.

14 \*Corresponding author(s). E-mail(s): [dak2173@columbia.edu](mailto:dak2173@columbia.edu);

15           Contributing authors: [anjali.das@columbia.edu](mailto:anjali.das@columbia.edu);

16           [clakhani@nygenome.org](mailto:clakhani@nygenome.org); [chloe.terwagne@crick.ac.uk](mailto:chloe.terwagne@crick.ac.uk);

17 [tlin@nygenome.org](mailto:tlin@nygenome.org); [tatsuhiko.naito@mssm.edu](mailto:tatsuhiko.naito@mssm.edu); [towfique.raj@mssm.edu](mailto:towfique.raj@mssm.edu);

18           **Abstract**

19           The increasing availability of whole-genome sequencing (WGS) has begun to elu-  
20 cide the contribution of rare variants (RVs), both coding and non-coding, to  
21 complex disease. Multiple RV association tests are available to study the relation-  
22 ship between genotype and phenotype, but most are restricted to per-gene models  
23 and do not fully leverage the availability of variant-level functional annotations.  
24 We propose Genome-wide Rare Variant EnRichment Evaluation (**gruyere**), a  
25 Bayesian probabilistic model that complements existing methods by learning  
26 global, trait-specific weights for functional annotations to improve variant prior-  
27 itization. We apply **gruyere** to WGS data from the Alzheimer's Disease (AD)  
28 Sequencing Project, consisting of 7,966 cases and 13,412 controls, to identify  
29 AD-associated genes and annotations. Growing evidence suggests that disruption  
30 of microglial regulation is a key contributor to AD risk, yet existing methods

31 have not had sufficient power to examine rare non-coding effects that incorpo-  
32 rate such cell-type specific information. To address this gap, we 1) use predicted  
33 enhancer and promoter regions in microglia and other potentially relevant cell  
34 types (oligodendrocytes, astrocytes, and neurons) to define per-gene non-coding  
35 RV test sets and 2) include cell-type specific variant effect predictions (VEPs)  
36 as functional annotations. `gruyere` identifies 15 significant genetic associations  
37 not detected by other RV methods and finds deep learning-based VEPs for splic-  
38 ing, transcription factor binding, and chromatin state are highly predictive of  
39 functional non-coding RVs. Our study establishes a novel and robust framework  
40 incorporating functional annotations, coding RVs, and cell-type associated non-  
41 coding RVs, to perform genome-wide association tests, uncovering AD-relevant  
42 genes and annotations.

43 **Keywords:** Rare variants, Alzheimer’s Disease, Bayesian probabilistic model,  
44 whole-genome sequencing

## 45 1 Main

46 The recent increase in available whole-genome sequencing (WGS) data has facilitated  
47 the study of rare variants (RVs), particularly in understanding their effects on complex  
48 diseases like late-onset Alzheimer’s disease (AD). AD is a neurodegenerative disorder  
49 with an estimated heritability between 59% and 74% [1]. While genome-wide associ-  
50 ation studies (GWAS) have identified over 100 loci linked to AD, with the *APOE-ε4*  
51 allele as the strongest genetic risk factor, they are restricted to common variant asso-  
52 ciations [2, 3]. Despite considerable efforts to quantify the polygenic nature of AD, a  
53 significant portion of genetic heritability remains unaccounted for. Some of this miss-  
54 ing heritability may be recovered with RVs [4]. RVs generally exhibit larger effect sizes  
55 than common variants, but their role is not yet well understood [5]. Studies have shown  
56 that integrating RVs into cumulative polygenic risk scores (PRS) can enhance predic-  
57 tive performance [6], but existing methods have identified fewer gene associations and  
58 have lower predictive power compared to common variant approaches. While a num-  
59 ber of genes, including *TREM2*, *ABCA7* and *SORL1* [7], have known RV associations  
60 in AD, the majority of these findings are restricted to coding variants. As most GWAS  
61 signals lie in the non-coding genome, expanding RV association studies beyond coding  
62 variants is critical. However, the study of non-coding RVs poses challenges due to the  
63 vast number of these variants, most of which likely have no functional impact [8]. It is  
64 therefore of substantial interest to use functional annotations for variant filtering and  
65 prioritization. To develop a more robust understanding of the contributions of both  
66 coding and non-coding RVs to AD, we propose a novel method that not only weights  
67 variants according to annotations but also prioritizes functional annotations that are  
68 most trait-relevant.

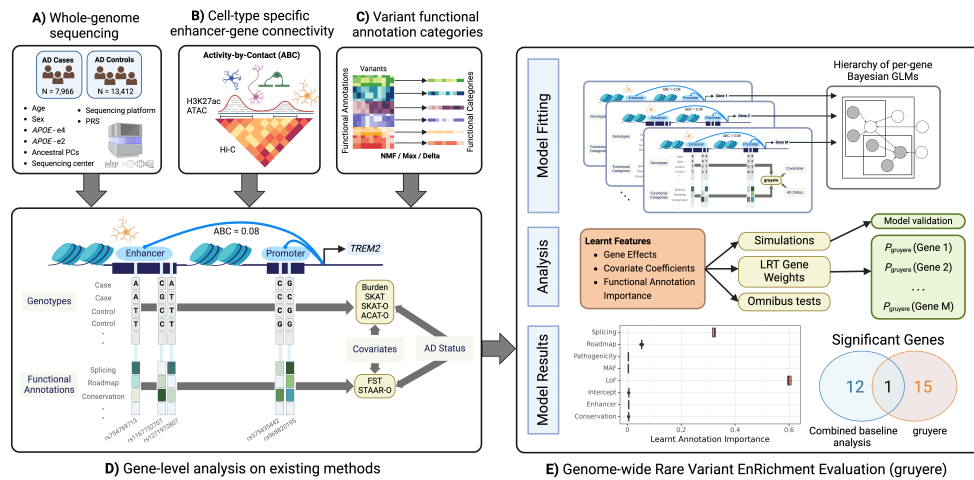
69 Applying traditional variant-level approaches like GWAS to RVs has low statistical  
70 power due to sparsity and a high multiple testing burden due to the large number of  
71 RVs compared to common variants. To address these limitations, RV methods aggre-  
72 gate variants in biologically related regions, typically by gene, to increase power [9].

73 More recent RV methods additionally account for functional annotations to prioritize  
74 relevant variants and filter out those predicted to have no function, which otherwise  
75 reduce power [10, 11]. Despite growing efforts to accurately predict which variants will  
76 affect particular molecular phenotypes (e.g., enhancer activation, RNA splicing) [12–  
77 16], there is a limited understanding of *which functions* are the most disease-relevant.  
78 Using functional annotations that have no phenotypic associations to weight RVs can  
79 add noise to models and decrease their power. This motivated us to develop a method  
80 that learns a genome-wide mapping from functional annotations to variant importance.

81 Growing evidence suggests that disrupted gene regulation in central nervous system  
82 (CNS) cell types, particularly microglia, is associated with the development and  
83 progression of AD [1, 17]. The majority of RV tests are developed for coding variant  
84 associations because 1) predicting functional coding variants is comparatively straight-  
85 forward (at least for loss-of-function), 2) population-scale whole-exome sequencing  
86 predates WGS, and 3) defining non-coding regions for testing is challenging in itself.  
87 Some methods use sliding windows, but testing overlapping windows of varying sizes  
88 can result in loss of power due to multiple testing [18]. Other methods use predicted  
89 *cis*-regulatory elements (CREs), in particular enhancers and promoters, to construct  
90 testing regions [19–21]. Given their modest size (typically less than 2kb), testing indi-  
91 vidual CREs still has limited statistical power. Combining multiple CREs that regulate  
92 a gene could help address this limitation but relies on accurate predictions of enhancer-  
93 gene links. We leverage the Activity-by-Contact (ABC) model, which predicts cell-type  
94 specific enhancer-gene connectivity using chromatin state and conformation data [22].  
95 We aggregate ABC-predicted enhancer-gene pairs to determine non-coding, cell-type  
96 and gene specific RV testing regions.

97 Due to the large number of genes and several million RVs found in population-scale  
98 WGS, existing methods are primarily restricted to per-gene models. This limits our  
99 understanding of disease-associated functional annotations. Most existing RV meth-  
100 ods are explicitly, or can be viewed as, generalized linear mixed models (GLMM). We  
101 instead develop a Bayesian generalized linear model (GLM) Genome-wide Rare Vari-  
102 ant EnRichment Evaluation, or **gruyere**, to model cell-type specific, non-coding RV  
103 associations on a genome-wide scale. In **gruyere**, a variant’s effect is a deterministic  
104 function of its annotations and the estimated AD-relevance of the gene it is linked  
105 to (if any). Our Bayesian model iteratively learns AD-relevant gene effects, covariate  
106 weights, and functional annotation importance while quantifying uncertainty, provid-  
107 ing increased flexibility to capture the complex, hierarchical structure of genetic data.  
108 We test our model using simulation analyses and compare results to several existing RV  
109 methods. We apply **gruyere** to WGS data from the Alzheimer’s Disease Sequencing  
110 Project (ADSP), consisting of 7,966 cases and 13,412 controls. Our model determines  
111 splicing, transcription factor (TF) binding, and chromatin state annotations most  
112 enriched for AD-associated non-coding RVs and identifies 16 significant genes, 15 of  
113 which are uniquely identified by **gruyere**. Of these, four – *C9orf78*, *MAF1*, *NUP93*,  
114 and *GALNT9* – remain significant in omnibus tests.

## 115 2 Results



**Fig. 1 Overview of the application of gruyere to AD.** Input data includes **A)** WGS and clinical information for AD cases and controls, **B)** Enhancer-gene interactions predicted by the ABC model for microglia, oligodendrocytes, astrocytes, and neurons, and **C)** variant functional annotations. **D)** Example analysis for the *TREM2* gene and microglia cell-type on existing methods. Columns represent RVs; light grey rectangles represent individual-level genotypes from WGS data for cases and controls; functional annotations for each RV are shown below genotypes; Burden, SKAT, SKAT-O, and ACAT-O are existing tests that use genotype, covariate, and AD status; FST and STAAR-O additionally use functional annotations. **E)** Workflow for gruyere. Per-gene RVs are aggregated and used for fitting the hierarchical Bayesian GLM. gruyere learns weights for covariates, genes, and functional annotations. We use simulations to assess gruyere at different heritabilities. Likelihood ratio tests are used to calculate gene-level *p*-values. Optionally, the gruyere *p*-values can be integrated with existing methods through omnibus testing.

### 116 2.1 Genome-wide Rare Variant EnRichment Evaluation 117 (gruyere) overview

**Table 1** Summary of gruyere Variables

Variable	Shape	Description
$Y$	$n \times 1$	Phenotypes for $n$ samples
$X$	$n \times c$	$c$ Covariates for $n$ samples
$G_g$	$n \times p_g$	Genotypes for $p$ variants in gene $g$ for $n$ samples
$Z_g$	$p_g \times q$	$q$ Functional annotations for $p$ variants in gene $g$
$\alpha_g$	$c \times 1$	$c$ Covariate coefficients
$\beta_{gj}$	$p_g \times 1$	Variant effect sizes for $p$ variants in gene $g$
$w_j$	$p_g \times 1$	Variant weight $\text{Beta}(MAF_j 1, 25)$
$w_g$	$1 \times 1$	Gene importance weights
$\tau$	$q \times 1$	Functional annotation weight

Current RV methods rely on independent per-gene models and, therefore, cannot capture genome-wide functional annotation importance. **gruyere** serves as a complementary method to existing RV tests by learning trait-specific functional annotation weights, covariate coefficients, and gene effects under a Bayesian framework (Supplemental Figure 1, Table 1). Rather than modeling each gene separately, we jointly fit **gruyere** as a hierarchy of per-gene GLMs using stochastic variational inference (SVI) [23]. We model AD risk for each gene  $g$ ,

$$\text{logit}(\mu_{ig}) = X_i\alpha_g + G_{ig}\beta_{gj}$$

where  $\mu_{ig}$  is the probability of AD for individual  $i$  given the genotypes for RVs associated with gene  $g$ ,  $X_i$  is a vector of covariates (e.g. sex, age, *APOE-e4* genotype), and  $G_{ig}$  is a genotype vector. We learn covariate weights  $\alpha_g$  and variant effects  $\beta_{gj}$ . We set  $\beta_{gj}$  to be a deterministic function of a learned gene effect  $w_g$ , transformed minor allele frequencies (MAFs)  $w_j$ , functional annotations  $Z$  (detailed in Methods 3.4 and Figure 1C), and learned annotation importance weights  $\tau$ ,

$$\beta_{gj} = w_g w_j (\tau_0 + \sum_{k=1}^q Z_{gjk} \tau_k).$$

118 In our analyses, **gruyere** learns annotation weights  $\tau$  for a range of annotations  $Z$ ,  
119 including in silico mutagenesis deep learning model predictions of splicing disruption  
120 (derived as the maximum of four individual SpliceAI scores [24]) and cell-type specific  
121 TF binding and chromatin state (derived from the Enformer model [14]). A larger  
122 magnitude of  $w_g$  indicates that disruption of gene  $g$  is associated with a higher pre-  
123 dicted risk of AD. Similar to a burden test, **gruyere** assumes all variants within a gene  
124 have the same direction of effect [25, 26]. However, because our functional annotations  
125 include both loss- and gain-of-function predictions, we are able to capture additional  
126 dispersion-based signal. To ensure robust generalization of learned parameters, we  
127 split data into training (80%) and test (20%) sets, where model weights are optimized  
128 using the training set and assessed on the unseen test set. We apply **gruyere** to both  
129 coding and non-coding RVs for AD, defining four cell-type specific non-coding groups  
130 for AD-relevant cell types (microglia, oligodendrocytes, astrocytes, and neurons [17])  
131 and testing each group individually.

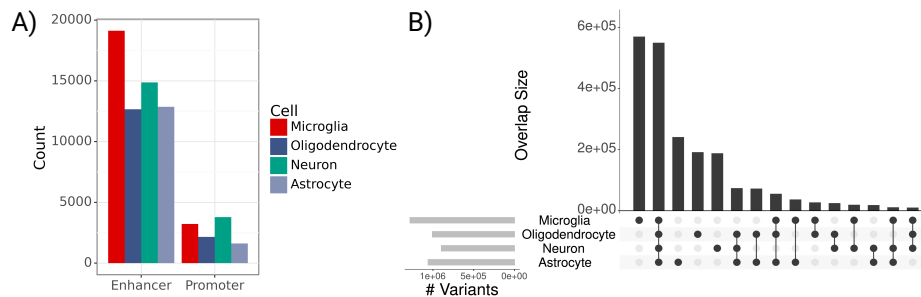
### 132 **Step 1. Estimating global annotation weights $\tau$ .**

133 Fitting  $\tau$  jointly across the entire genome would be 1) computationally challenging  
134 due to the large number of RVs and 2) statistically inefficient, as only AD-associated  
135 genes will contribute relevant signal. We therefore estimate  $\tau$  under the **gruyere** model  
136 from a subset of potentially AD-relevant genes identified using a lenient significance  
137 threshold (nominal  $p < 0.01$ ) for the Functional Score Test (FST)[27]. We assess  
138 the robustness of **gruyere** estimates when selecting genes with varied significance  
139 thresholds and for a number of existing RV tests and find annotation weights  $\tau$  are  
140 broadly consistent ( $+/- 0.02$ ).

141 **Step 2. Per-gene analysis.**

142 Once genome-wide estimates for  $\tau$  are obtained, **gruyere** simplifies to a logistic  
 143 regression that learns covariate  $\alpha_g$  and gene  $w_g$  weights. We efficiently fit **gruyere**  
 144 separately and in parallel for all genes, holding  $\tau$  fixed. We perform likelihood ratio  
 145 tests (LRT) to compare a covariate-only regression against combined covariate and  
 146 genotype regression models to determine gene-level significance for  $w_g$ .

147 **2.2 Constructing cell-type & gene specific variant sets using**  
 148 **predicted CREs**



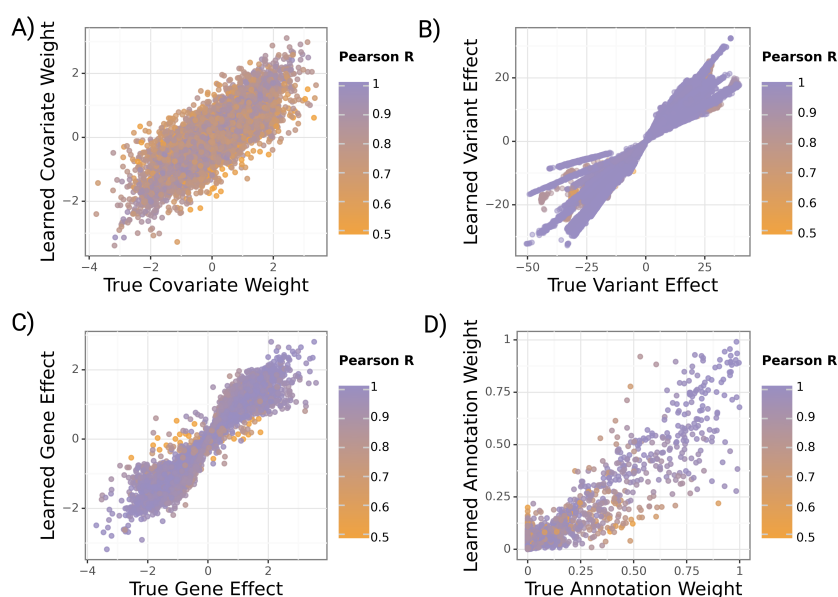
**Fig. 2 Predicted regulatory element and variant counts across cell types. A)** Bar plot of predicted CRE counts by cell type ( $ABC > 0.02$ ). **B)** Upset plot of variant overlap across 4 cell types in ADSP data; Light grey bars on left indicate total RV counts for each cell type; Vertically connected dots represents groups and corresponding bars indicate variant overlap for that group.

149 We grouped non-coding RVs by gene and CNS cell type using the Activity-by-  
 150 Contact (ABC) model (Figure 1B) [22]. The ABC model uses epigenomic profiles and  
 151 chromatin conformation to determine cell-type specific enhancer-gene interactions,  
 152 filtering out genes that are not expressed. We use publicly available ATAC-seq and  
 153 H3K27ac ChIP-seq signals for microglia, oligodendrocytes, astrocytes, and neurons  
 154 [28], as well as Hi-C averaged across ten CNS cell types to account for 3-dimensional  
 155 chromatin interactions. For each gene, we analyze RVs aggregated across all CREs  
 156 interacting with that gene ( $ABC > 0.02$ ). We test each cell type separately and also  
 157 analyze rare coding variants for comparison. In total, ABC defines 70,300 CREs across  
 158 all four cell types and 17,929 genes, with higher relative counts of microglia-predicted  
 159 enhancers (Figure 2A). Predicted CREs frequently co-occur across cell types, with  
 160 39.4% of CREs found in more than one cell type. Promoter regions tend to have higher  
 161 ABC scores than enhancers (mean  $ABC = 0.07$  vs. mean  $ABC = 0.04$ ), but their  
 162 genomic lengths are similar, with an average length of 632bp and standard deviation  
 163 of 132bp. ABC accounts for interactions of a single enhancer with multiple genes so  
 164 one RV can be linked to multiple genes. In our analysis, an ABC enhancer maps to  
 165 an average of 3.8 genes in microglia and between 5.4 and 5.9 genes in the other three  
 166 cell types. Our non-coding variant sets contain an average of 376 RVs per gene. There  
 167 are a total of 2,092,931 RVs in CREs across the four cell types, 901,570 of which are

168 included for more than one cell type, and 550,001 that are in all four cell types (Figure  
169 2B).

### 170 2.3 Simulation studies confirm accurate estimation of model 171 parameters

172 We generate synthetic phenotypes (see Methods 3.5) and fit `gruyere` on 100 sets of  
173 simulated data with estimated heritability between 5% and 30% (detailed in Supple-  
174 mental Methods [29]) using 500 randomly selected genes. We find that all variables  
175 are well recovered, with average Pearson correlations  $R = 0.81, 0.95, 0.98, 0.97$  for  
176  $\alpha_g, \beta_{gj}, w_g$  and  $\tau$  respectively (Figure 3). Covariates  $\alpha_g$  have the lowest  $R$ , possibly  
177 due to correlated covariates. Average recovery across all variables remains high when  
178 varying the prior distributions (Pearson  $R > 0.78$ ) as well as when simulated distri-  
179 butions differ from the priors used during inference (Pearson  $R > 0.66$ ). Results are  
180 robust to the number of covariates, genes, and annotations modeled.

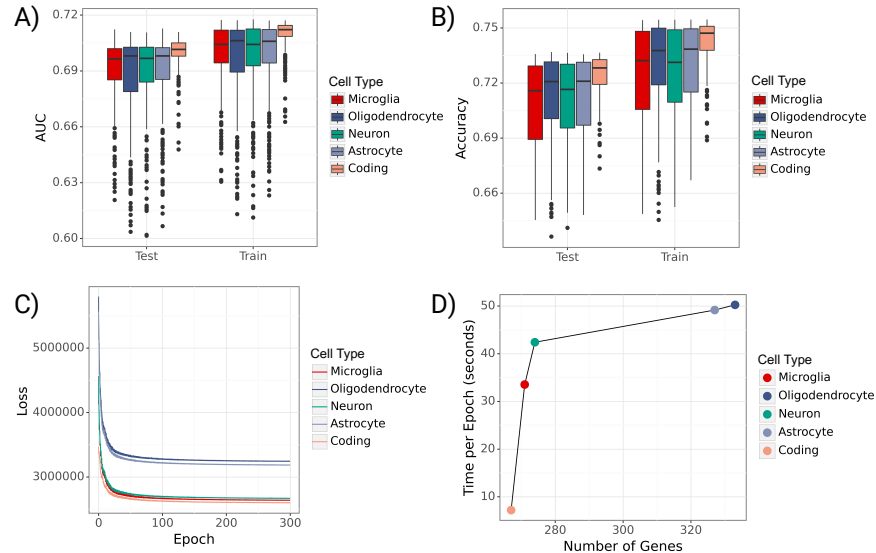


**Fig. 3** Learned versus true `gruyere` parameters across 100 simulations. Points are colored by the Pearson correlation coefficient of a parameter for a given simulation. **A)** Covariate regression coefficients  $\alpha$  ( $c = 30$  covariates). **B)** Variant effect  $\beta$ . **C)** Gene effect  $w_g$  ( $M = 500$  randomly selected genes). **D)** Annotation weight  $\tau$  ( $q = 13$  annotations).

181 We analyze how simulation performance correlates with overall and genetic her-  
182 itability for each simulation. This allows us to more meaningfully evaluate model  
183 performance for complex diseases like AD where estimated heritability is low. As  
184 expected, we find that `gruyere` is better able to recover  $\beta_{gj}$  and  $w_g$  with increased  
185 genetic heritability (Supplemental Figure 2). However, even when total heritability

186 is as low as 5%, the minimum correlation between true and estimated parameters  
 187 remains quite high (Pearson  $R = 0.68$ ).

## 188 2.4 Applying gruyere to AD WGS data reveals novel disease 189 associations

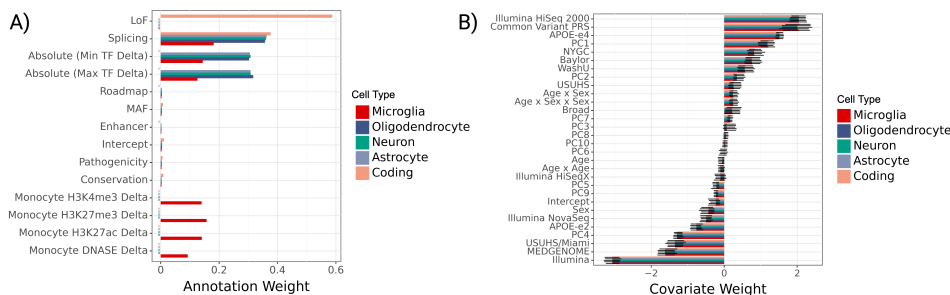


**Fig. 4 Performance of gruyere.** **A)** Boxplots of per-gene AUROCs for train and test sets across cell types. **B)** Boxplots of per-gene accuracies for train and test sets across cell types. **C)** Trace ELBO loss over 300 epochs across cell types. **D)** Average training time per epoch (seconds) versus number of genes used in joint model for each cell type.

190 **Performance on ADSP WGS data.** After validating model performance through  
 191 simulations, we fit *gruyere* to the ADSP WGS data. We analyze coding and non-  
 192 coding (microglia, oligodendrocyte, neuron, astrocyte) groups separately, and refer to  
 193 each set as a cell type. For each cell type, we use a subset of genes for joint fitting  
 194 (FST  $p$ -value  $< 0.01$ ), leading to between 267 and 333 genes per cell-type. AD predic-  
 195 tion performance is fairly consistent across non-coding variants (average test set Area  
 196 Under the Receiver Operating Characteristic, or AUROC, of 0.69) and slightly higher  
 197 for coding variants (average  $AUROC_{test} = 0.70$ ) (Figure 4A-B). When averaging predic-  
 198 ted probabilities across genes, performance further improves ( $AUROC_{test} = 0.72$ )  
 199 for all cell types. These metrics are in line with current AD literature and outperform  
 200 a covariate-only regression model ( $AUROC_{test} = 0.65$ ) [30]. There is a substantial  
 201 range in prediction performance for each cell type (e.g. minimum  $AUROC_{test} = 0.62$ ,  
 202 maximum  $AUROC_{test} = 0.71$  for microglia), highlighting the varying degrees of asso-  
 203 ciation with AD across genes. We find that gene-level performance is consistent across  
 204 model refitting and that the loss converges reliably (4C) [31]. Fitting time increases



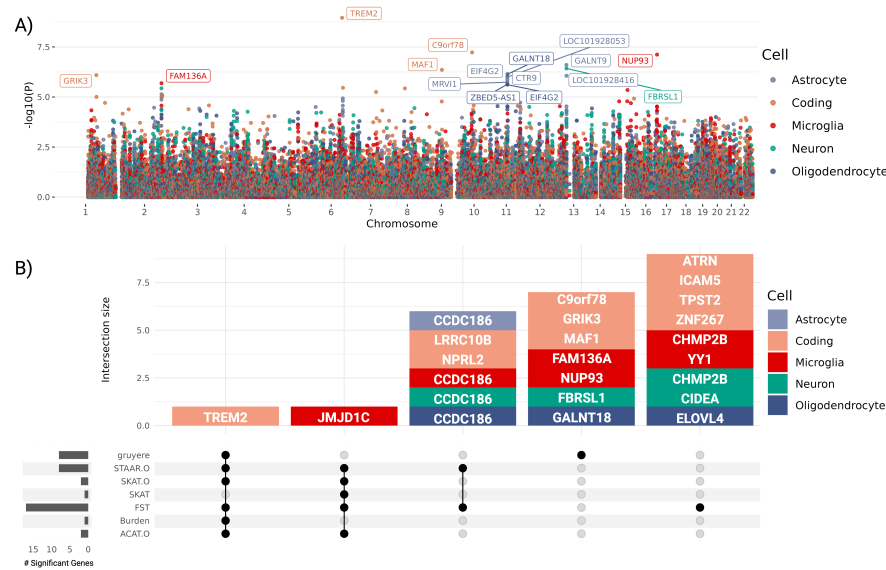
205 approximately logarithmically with the total number of genes (Figure 4D). On average,  
 206 it takes 37 seconds per epoch and three hours total to jointly fit `gruyere` across  
 207 300 epochs. Per-gene estimation is much faster, taking an average of 4.3 seconds per  
 208 gene to complete.



**Fig. 5** `gruyere` parameters learned from ADSP WGS. **A)** Bar plot of genome-wide annotation weights  $\tau$  learned in jointly fit model across cell types. We denote crosses (X) to the left of bars where an annotation is not included for a cell type. **B)** Bar plot of per-gene covariate weights ( $\alpha_g$ ) learned in jointly fit model across cell types. Error bars illustrate the minimum and maximum values learned across genes.

209 **Learned annotation and covariate weights.** We find that the top `gruyere`  
 210 functional annotation weights come from splicing across all non-coding RV groups and  
 211 loss-of-function (LoF) for coding variants (Figure 5A). LoF variants can be highly  
 212 disruptive to gene function and are often used as a variant filtering method in gene-  
 213 based tests. Therefore, it is predictable that we find `gruyere` places a large weight  
 214 on LoF coding RVs. It is perhaps not surprising that `gruyere` also prioritizes RVs  
 215 predicted to disrupt normal splicing, as they can substantially change the protein  
 216 product or have large effects on gene dosage via nonsense mediated decay. For all  
 217 non-coding regions, cell-type specific TF binding predictions from Enformer contain  
 218 the next largest annotation weights. This suggests RVs associated with an increase  
 219 (Max TF Delta) or decrease (Min TF Delta) in binding are predicted to have larger  
 220 effects on AD, at least in AD-relevant genes. For microglia RV sets, we additionally  
 221 find increased AD association for variants related to histone modification (H3K4me3,  
 222 H3K27me3, H3K27ac) and DNASE in monocytes (often used as a proxy for microglia  
 223 [32]). We restrict Enformer annotations to non-coding variants as they are specific  
 224 to cell-types. The enhancer category has very small weights across cell-types. Since  
 225 all variants included in the non-coding analyses are in putative CREs, it is perhaps  
 226 not too surprising that cell-type specific enhancer annotations are lowly prioritized by  
 227 `gruyere`.

228 Covariate effects are learned consistently across genes and cell types, with sequenc-  
 229 ing center, common variant polygenic risk score, and Illumina HiSeq 2000 sequencing  
 230 platform as the top three covariates (Figure 5B). As expected, *APOE-e4* is learned to  
 231 have a large positive risk effect while the *APOE-e2* allele has a negative (protective)  
 232 effect [2]. These effects agree well with those of a simple logistic regression predicting  
 233 AD status from covariates.



**Fig. 6 Top gruyere genes. A)** Manhattan plot across cell types. The Y-axis shows  $-\log_{10}(p\text{-value})$  for each gene and X-axis shows gene position. Each color is a cell type. **gruyere**-significant genes are labeled. **B)** Stacked upset plot of significant gene overlap across all tests after pruning for coregulation. Dark grey bars on left indicate total number of significant genes for each test. Vertically connected dots represent groups and corresponding bars indicate the number of overlapping significant genes identified for that group. Each bar is colored by cell type.

234 **Learned gene effects and associations:** Estimated gene effects  $w_g$  are fit from  
 235 a per-gene logistic regression, where we use LRTs to determine gene-level **gruyere**  $p$ -  
 236 values (Methods 3.2). Significant genes after Bonferroni correction for each cell type  
 237 are shown in Figure 6A, where a total of 16 genes reach genome-wide significance.  
 238 The well-established *TREM2* RV association [33], as well as *MAF1*, *C9orf78*, and  
 239 *GRIK3* are found significant for coding variants. Although not as widely recognized as  
 240 *TREM2*, *MAF1* has been previously reported in association with AD [34], and *C9orf78*  
 241 has been identified in an AD dementia meta-analysis [35]. *GRIK3* has emerged as a  
 242 gene of interest due to the role of kainate receptors in neuroinflammation, a key feature  
 243 of AD. Inflammatory responses can amplify glutamate release and disrupt receptor  
 244 functioning, which may further accelerate neurodegeneration. This makes *GRIK3*,  
 245 and glutamate signaling more broadly, potential targets for therapies [36, 37]. The  
 246 identification of these genes by **gruyere** highlights their potential as candidate genes  
 247 for further study in AD.

248 **gruyere** identified 12 non-coding RV associations across cell types with 2 in  
 249 microglia, 6 in astrocytes, 1 in neurons, and 3 in oligodendrocytes. The most significant  
 250 of these is *NUP93* (microglia), which, although not specifically linked to AD, is part of  
 251 a group of nucleoporin (Nup) mutations associated with neurodegenerative disorders  
 252 like AD [38]. Three significant genes, *GALNT9*, *FBRSL1*, and *LOC101928416*, are  
 253 closely located on chromosome 12q24.33 and share over 80% of their ABC-predicted  
 254 CREs, indicating that their associations are driven by the same set of RVs. Although

255 variants can map to multiple genes in our model framework, making our analysis sus-  
256 ceptible to coregulation, we are able to investigate and identify the specific CREs  
257 driving these associations. Of the overlapping promoters for these three genes, regions  
258 have higher ABC scores for *FBRSL1* (ABC = 0.20) compared to *GALNT9* (ABC =  
259 0.04) and *LOC101928416* (ABC = 0.06), suggesting a stronger regulatory impact on  
260 *FBRSL1* and further isolating overlapping signal. *FBRSL1* (neuron), has not been  
261 linked to AD, but it presents a strong candidate gene for its distinctive neuronal expres-  
262 sion profile and involvement in neurogenesis and transcriptional regulatory networks  
263 [39]. Multiple associations (*GALNT18*, *CTR9*, *EIF4G2*, *ZBED5-AS1*, *LOC101928053*,  
264 and *MRV11*) specific to astrocyte and oligodendrocyte cell-types are coregulated in  
265 chromosome 11p15.4, and the strongest signal, *GALNT18*, has been connected to AD  
266 in more than one study [40, 41]. After pruning coregulated signals, **gruyere** identifies  
267 8 significant genes.

268 We identify both known and novel AD-associated risk genes with **gruyere**. Sig-  
269 nificant **gruyere** genes are associated with increased gene expression across thirteen  
270 brain tissues found in GTEx (2-sample t-test  $p = 9.2 \times 10^{-31}$ , Supplemental Figure  
271 3A) [42]. 5 of our 16 significant genes have expression QTLs in our microglia genomic  
272 atlas (isoMiGA) that colocalize with AD or Parkinson’s disease (PD) GWAS (Sup-  
273 plemental Figure 3B) [43–46]. Specifically, *TREM2* and *MAF1* have significant SNPs  
274 in a recent AD GWAS [44] while *FBRSL1*, *EIF4G2*, and *ZBED5-AS1* are signifi-  
275 cant in a large PD GWAS [45]. AD and PD have known genetic overlap, motivating  
276 QTL colocalization of both traits [47]. Finally, we compare **gruyere**  $p$ -values with the  
277 Alzheimer’s Disease Variant Portal (ADVP) catalog of 956 reported AD genes, find-  
278 ing that **gruyere** yields more significant  $p$ -values for ADVP versus non-ADVP genes  
279 (2-sample t-test  $p = 7.1 \times 10^{-9}$ , Supplemental Figure 3C) [48].

#### 280 *Comparing gruyere to existing methods.*

281 We compare pruned **gruyere** results with AD associations identified by a number of  
282 existing RV methods: burden test, sequence kernel association test (SKAT), optimal  
283 unified test (SKAT-O), functional score test (FST), aggregated Cauchy association test  
284 (ACAT), and variant-set test for association using annotation information (STAAR)  
285 [10, 25, 27, 49–51]. Burden, SKAT, SKAT-O, and ACAT-O tests do not include  
286 functional annotations, while FST and STAAR incorporate them (description in Sup-  
287 plemental Table 1 and detailed in Figure 1D for *TREM2* and microglia). We use the  
288 same set of functional annotations for FST and STAAR as for **gruyere**. We find that  
289 **gruyere**  $-\log_{10}(p\text{-values})$  have the highest correlation with Burden tests (Pearson  
290  $R = 0.86$ ) and show moderate to high correlation with combination methods STAAR-  
291 O, SKAT-O, FST, and ACAT-O (Pearson  $R = 0.45 - 0.58$ ) (Supplemental Figure  
292 4A). The higher observed correlation with burden tests is expected, as **gruyere** also  
293 assumes unidirectional variant effects within a gene. We examine overlap of significant  
294 genes across all tests and find that there is minimal overlapping signal across meth-  
295 ods (Figure 6B). Of the 16 (8 pruned) AD associations identified by **gruyere**, 15 (7  
296 pruned) are unique to **gruyere**, while *TREM2* (coding) is detected across all tests but  
297 SKAT where it narrowly misses significance. In total, burden, SKAT, SKAT-O, and

298 ACAT-O tests identify only two significant associations, highlighting the importance  
299 of including functional annotations, particularly for non-coding RV associations.

300 **Integrating gruyere into omnibus tests.** We combine *gruyere*  $p$ -values with  
301 existing methods using ACAT (Supplemental Figure 4B). Comparing *ACAT*(burden,  
302 SKAT, SKAT-O, FST, ACAT, STAAR) to *ACAT*(*gruyere*, burden, SKAT, SKAT-  
303 O, FST, ACAT, STAAR), we find that the inclusion of *gruyere* in omnibus tests  
304 boosts the number of significant associations from 12 to 16, adding *C9orf78*, *MAF1*,  
305 *NUP93* and *GALNT9*. There is no loss of power with this method, as all existing  
306 signals remain after including *gruyere*; we simply increase the total number of AD  
307 associations identified.

## 308 3 Methods

### 309 3.1 Data Overview

310 **Whole-Genome Sequencing Data:** We analyze the latest release of WGS data  
311 from the Alzheimer’s disease sequencing project (ADSP) consisting of 21,378 unre-  
312 lated individuals over the age of 65 (7,966 cases, 13,412 controls) after QC [52? ]. We  
313 follow a standard pipeline to QC WGS data. First, we combine phenotype information  
314 across multiple cohorts and remove genetically identical duplicates ( $IBD \hat{\pi} > 0.95$ ) and  
315 technical replicate samples, selecting samples with the highest call rates. We priori-  
316 tize phenotype information for individuals in family studies over case-control studies.  
317 Related individuals are removed using Kinship-based Inference for Gwas (KING) [53],  
318 keeping AD cases where possible. In PLINK [54], we remove individuals with more  
319 than 10% genotype missingness, variants with less than 90% genotyping rate, and  
320 keep only biallelic variants with an observed  $MAF \leq 0.05$ . Missing genotypes are  
321 imputed as the average observed MAF. For analysis, we randomly split samples into  
322 80% train and 20% unseen test sets, stratifying by ancestry. ADSP samples are pri-  
323 marily of European ( $N = 9,133$ ), African ( $N = 5,173$ ) and Hispanic (5,059) ancestry,  
324 with smaller South Asian ( $N = 1,951$ ) and East Asian ( $N = 62$ ) groups.

325 **Clinical Information:** We use 30 available covariates in our model: sex, age,  
326  $age^2$ ,  $age \times sex$ ,  $age \times sex^2$ , *APOE-e4* genotype, *APOE-e2* genotype, 10 ancestry  
327 principal components calculated from the 1000 Genomes Project [55], a common vari-  
328 ant PRS [56], one-hot encoded sequencing platform (Illumina HiSeq 2000, HiSeqX,  
329 Nova Seq), one-hot encoded sequencing center (Illumina, USUHS, USUHS/Miami,  
330 NYGC, MEDGENOME, Baylor, Broad, WashU), and an intercept term (Figure 1A).  
331 Covariates are min-max scaled to a range of 0 to 1.

### 332 3.2 Proposed Bayesian rare variant model: gruyere

We develop *gruyere*, a hierarchy of per-gene GLMs (Supplemental Figure 1). We  
define our model jointly as

$$\text{logit}(\mu_{ig}) = X_i \alpha_g + G_{ig} \beta_{gj}$$

where  $\mu_{ig}$  is the probability of AD for individual  $i$  associated with gene  $g$ ,  $X_i$  is a vector of covariates, and  $G_{ig}$  is a vector of genotype dosages for each RV. Covariate coefficients  $\alpha_g$  are modeled from prior,

$$\alpha_g \sim \text{Normal}(0, 1).$$

The key innovation in the **gruyere** model is the construction of per-variant genetic effects for gene  $g$ , or  $\beta_{gj} = (\beta_{gj1}, \dots, \beta_{gjp})^T$ , which is defined as the product of gene effects, transformed MAF, and weighted functional annotations. Of note, if a variant  $j$  is included in the RV set for both genes  $g_1$  and  $g_2$ , the variant effect can differ for  $\beta_{g_1j}$  and  $\beta_{g_2j}$ . We define  $\beta_{gj}$  as

$$\beta_{gj} = w_g w_j (\tau_0 + \sum_{k=1}^q Z_{gjk} \tau_k)$$

where  $w_g$  are gene importance weights,  $w_j$  are variant weights based on observed MAF as suggested by Wu et al. [49],  $\tau$  are genome-wide annotation importance scores, and  $Z_{gjk}$  is a scaled functional annotation  $k$  for RV  $j$  and gene  $g$ . The variables are modeled as,

$$\begin{aligned} w_j &= \text{Beta}(MAF_j | 1, 25) \\ w_g &\sim \text{Normal}(0, 1) \\ \tau &\sim \text{Dirichlet}(\mathbf{1}_q) \end{aligned}$$

333 For each gene, we use a Bernoulli likelihood to sample  $\sigma(X_i \alpha_g + G_{ig} \beta_{gj})$ , and  
334 aggregate loss across each  $g \in M$ . Learned parameters are  $\alpha_g, w_g$ , and  $\tau$ . We select  
335 a Dirichlet prior for annotation weights  $\tau$  to ensure identifiability between  $\tau$  and  $w_g$ .  
336 Without constraining  $\tau$  to a fixed sum,  $w_g$  can be swapped for  $w_g/c$  and  $\tau$  for  $c\tau$  for  
337 any positive constant  $c$  without changing the likelihood, leading to non-identifiability  
338 between gene and annotation weights.

339 We approximate the true posterior distribution for **gruyere** by minimizing the  
340 Kullback-Leibler (KL) divergence, which is equivalent to maximizing the Evidence  
341 Lower Bound (ELBO) [31]. To maximize the ELBO, we use SVI, implemented in the  
342 `pyro` probabilistic programming language [57]. We approximate the posterior distribu-  
343 tion of latent variables  $\alpha_g$  and  $w_g$  with mean field normal distributions (AutoNormal  
344 guide), while optimizing annotation weights  $\tau$  as point estimates with a Delta distribu-  
345 tion (AutoDelta guide). We apply the Adam optimizer, a learning rate of 0.1, train  
346 for 300 epochs, and draw 50 samples from the posterior to estimate standard devi-  
347 ations of the learned parameters. We explore different prior distributions for all key  
348 parameters.

349 Once global  $\tau$  is learned, we streamline **gruyere** with a per-gene analysis. Holding  $\tau$   
350 fixed, our model simplifies to a logistic regression where only  $\alpha_g$  and  $w_g$  are estimated.  
351 **gruyere** efficiently computes gene-level  $p$ -values using LRTs comparing a covariate-  
352 only regression to a combined covariate and genotype regression model:

$$LR_g = -2 \times (LL_{\text{combined}_g} - LL_{\text{covariate only}}), \quad df = 1$$

$$\text{gruyere } p\text{-value}(g) = P(\chi^2 > LR_g)$$

where  $LR$  are the log-likelihoods for each logistic regression. For each cell type, we use Bonferroni correction to define the  $p$ -value significance threshold:

$$p < \frac{0.05}{\# \text{ genes per cell type}} \Rightarrow 2.88 \times 10^{-6} < p < 3.64 \times 10^{-6}$$

### 3.3 Cell-type specific RV gene set prediction using the ABC model

We calculate enhancer-gene connectivity using publicly available ATAC-seq and H3K27ac ChIP-seq data for human microglia, oligodendrocytes, astrocytes, and neurons [28]. We apply ABC to this data following the guidelines and default parameters provided at <https://github.com/broadinstitute/ABC-Enhancer-Gene-Prediction>. This involves first calling candidate peak regions for the ATAC-seq DNase hypersensitive sites (DHS) using MACS2 (peakExtendFromSummit = 250, nStrongestPeaks = 150000). Then we quantify enhancer activity as the geometric mean of the read counts of DHS and H3K27ac ChIP-seq in candidate enhancer regions. Lastly, we compute the ABC score using averaged Hi-C data (hic\_resolution = 5000) fit to the power-law model. The omics data is aligned to hg19, so we converted the ABC-predicted start and end positions of enhancers to hg38 for analysis. For each gene and separately each cell type, we aggregate all elements  $E$  for gene  $G$  that have an  $ABC \geq 0.02$  and extract RVs from within these regions to determine our cell-type specific non-coding RV gene sets.

### 3.4 Calculating functional annotation groups

We use a range of variant-level functional annotations primarily from the Whole Genome Sequencing Annotation database [58]. Annotations with greater than 5% missingness in our RVs are removed, resulting in 50 coding variant and 52 non-coding variant functional annotations listed in Supplemental Table 2. To reduce dimensionality of related annotations while accounting for their diverse biological effects, we apply non-negative matrix factorization (NMF) to summarize groups of related annotations, inspired by STAAR [10]. We use NMF to retain interpretable directional scaling of annotations. Based on correlation structure and a priori knowledge, we define six major functional categories – splicing, conservation, integrative deleterious predictions, brain-related Roadmap epigenetics, population-specific MAF, and enhancer activity [15, 59, 60]. Because the splicing annotation group is derived from four SpliceAI predictions that are not highly correlated and sparsely distributed, we instead use the maximum score for this category. We also include a binary LoF prediction calculated with Loss-Of-Function Transcript Effect Estimator (LOFTEE) [61] for coding variants along with an intercept term. All annotations are scaled between 0 and 1, where a larger value represents increased predicted variant function.

**Deep Learning Delta Scores:** For all four cell types, we include additional cell-type specific functional annotations: absolute maximum and absolute minimum TF delta scores derived from Enformer [14], a deep learning genomics model. We calculate

389 variant delta scores for 5,318 functional genomics assays. The Enformer model predicts  
390 read counts (in 128 BP bins) of these assays as a function of 196,608 BP input DNA  
391 sequence. For a particular variant, we compare the model output of the reference  
392 sequence, centered around the variant position, with the output of the alternative  
393 sequence which replaces the reference allele with the alternative allele. For a particular  
394 genomics assay, the delta score is the difference between the sum of reference sequence  
395 predictions for the middle 32 bins and the sum of the alternate sequence predictions  
396 for the same bins. We normalize these scores by first calculating the delta scores for  
397 the approximately 18 million variants from the UK Biobank cohort used in PolyFun  
398 [62, 63], and then Z-score normalize each assay according based on this collection of  
399 variants. We apply this normalization to the delta scores used in our analysis. We  
400 aggregate delta scores to determine composite maximum and minimum predictions  
401 for each variant, highlighting the delta scores of only the enriched TFs within each of  
402 the four cell types ([28]). For microglia non-coding variant sets, we additionally use  
403 delta scores for 4 epigenetic marks (H3K4me3, H3K27ac, H3K27me3, and DNASE)  
404 for monocytes, a proxy for microglia.

### 405 3.5 Data Simulation

We use simulations to evaluate **gruyere** performance. We randomly sample values for each parameter and use these simulated variables in the GLM framework,  $\text{logit}(\mu_{ig}) = X_i\alpha_g + G_{ig}\beta_g$ . The real ADSP genotypes  $G_{ig}$ , covariates  $X_i$ , and functional annotations  $Z_{gjk}$  along with simulated parameters  $\alpha_g^S, w_g^S$ , and  $\tau^S$ , generate simulated phenotypes  $Y_{ig}^S$ . Simulations are restricted to a maximum estimated heritability of 30% to realistically evaluate complex diseases. For each simulation, we draw **gruyere** parameters from the following distributions:

$$\alpha_g^S \sim \text{Normal}(0, 1), \quad w_g^S \sim \text{Normal}(0, 1), \quad \tau^S \sim \text{Dirichlet}(\mathbf{1}_q)$$

We define  $\beta_{gj}$  in the same way, simply using simulated variables:

$$\beta_{gj}^S = w_g^S w_j \sum_{k=1}^q Z_{gjk} \tau_k^S$$

406 Using this combination of true data and simulated parameters, we sample synthetic  
407 phenotypes  $y_g^S$  from a Bernoulli distribution. We perform 100 simulations on  $M =$   
408 500 randomly selected genes. In general, we sample **gruyere** parameters from the  
409 same distribution that they are learned. We have also tested model performance when  
410 simulated data comes from a different distribution than its learned counterpart (e.g.  
411 sampling  $w_g$  from a Normal distribution in simulations but fitting from a Gamma  
412 prior).

## 413 4 Discussion

414 We develop **gruyere**, a functionally-informed RV association test that fits a hierar-  
415 chy of Bayesian GLMs to estimate genome-wide functional annotation importance,

416 gene effects, and covariate coefficients. **gruyere** builds upon existing RV methods with  
417 two key advancements: 1) a genome-wide approach that enables trait-specific weight-  
418 ing of functional annotations, and 2) a flexible, powerful and calibrated probabilistic  
419 framework that estimates uncertainty. We incorporate an innovative methodology for  
420 analyzing RVs in the non-coding genome. Using the Activity-by-Contact model, we  
421 predict cell-type-specific enhancer-gene connectivity from chromatin state and con-  
422 formation data, aggregating predictions by gene to define interpretable non-coding  
423 RV testing regions. We use in silico mutagenesis under state-of-the-art deep learning  
424 models of pre (SpliceAI) and post (Enformer) transcription gene regulation to pre-  
425 dict RV effects. Simulation analyses validate **gruyere** and show it is able to recover  
426 ground truth parameters across diverse model specifications and even for realistically  
427 low heritability.

428 We apply **gruyere**, along with a number of established RV association tools, to  
429 the most recent WGS release from ADSP. Our analysis identifies both known (e.g.,  
430 *TREM2*) and novel (e.g., *NUP93*) candidate AD genes. Specifically, **gruyere** uniquely  
431 identifies 15 genes, of which *C9orf78*, *MAF1*, *NUP93* and *GALNT9* remain signif-  
432 icant in aggregated Cauchy tests. Our analysis additionally provides an improved  
433 understanding of AD-relevant functional annotations. **gruyere** confirms the expecta-  
434 tion that LoF is the most informative annotation for coding variants, but additionally  
435 finds deep learning-based predictions for splicing, TF binding and chromatin state are  
436 highly predictive of functional non-coding RVs.

437 We use ancestry principal components as covariates to account for population  
438 diversity, but one area for future work would be integrating a random effect term to  
439 better account for relatedness and population structure [64]. Another possible exten-  
440 sion to **gruyere** would be incorporating gene-level features as priors [65]. While we  
441 focus our analysis on AD, **gruyere** can be applied to any complex disease with suf-  
442 ficient WGS data. As the quality of functional annotations continues to improve,  
443 **gruyere** will become an increasingly valuable tool for identifying disease-associated  
444 genes and annotations.



445 **Supplementary information.** Supplementary Figures 1-4; Supplementary Table  
446 1-2; Supplementary Methods; Supplementary Acknowledgements; Supplementary  
447 Data: Gene  $p$ -values for each cell-type.

448 **Acknowledgements.** Research reported in this paper was supported by  
449 Alzheimer's Disease Sequencing Project of the National Institutes of Health under  
450 award number 5 U01 AG068880-02. The content is solely the responsibility of the  
451 authors and does not necessarily represent the official views of the National Institutes  
452 of Health. We are thankful to Tulsi Patel, Edoardo Marcora, Alison Goate, and Iuliana  
453 Ionita-Laza for helpful feedback and discussions.

## 454 **Declarations**

- 455 • Funding: Research reported in this paper was supported by Alzheimer's Disease  
456 Sequencing Project of the National Institutes of Health under award number 5 U01  
457 AG068880-02.
- 458 • Declaration of interests: T.R. served as a scientific advisor for Merck and serves as  
459 a consultant for Curie.Bio.
- 460 • Ethics approval and consent to participate: Not applicable
- 461 • Consent for publication: Not applicable
- 462 • Data availability: This paper uses the ADSP Release 4 WGS data and AD phenotype  
463 data.
- 464 • Materials availability: Not applicable
- 465 • Code availability: Code and details for running `gruyere` is available on GitHub:  
466 <https://github.com/daklab/gruyere>
- 467 • Author contribution: D.A.K. conceived and supervised the project. A.D. and D.A.K.  
468 developed the methods. A.D. and D.A.K. wrote the manuscript. A.D. wrote the  
469 software code and performed the analyses. All authors read, reviewed, and approved  
470 the final manuscript.

471 Editorial Policies for:

472 Springer journals and proceedings: <https://www.springer.com/gp/editorial-policies>

473 Nature Portfolio journals: <https://www.nature.com/nature-research/editorial-policies>

474 *Scientific Reports*: <https://www.nature.com/srep/journal-policies/editorial-policies>

475 BMC journals: <https://www.biomedcentral.com/getpublished/editorial-policies>

## 476 **References**

- 477 [1] Sims, R., Hill, M., Williams, J.: The multiplex model of the genetics of  
478 alzheimer's disease. *Nat. Neurosci.* **23**, 311–322 (2020) <https://doi.org/10.1038/s41593-020-0599-5>  
479

- 480 [2] Andrews, S.J., Renton, A.E., Fulton-Howard, B., Podlesny-Drabiniok, A., Mar-  
481 cora, E., Goate, A.M.: The complex genetic architecture of alzheimer's disease:  
482 novel insights and future directions **90**, 104511 (2023) [https://doi.org/10.1016/j.](https://doi.org/10.1016/j.ebiom.2023.104511)  
483 [ebiom.2023.104511](https://doi.org/10.1016/j.ebiom.2023.104511)
- 484 [3] Auer, P.L., Lettre, G.: Rare variant association studies: considerations, challenges  
485 and opportunities **7**, 16 (2015) <https://doi.org/10.1186/s13073-015-0138-2>
- 486 [4] Wainschtein, P., Jain, D., Zheng, Z., TOPMed Anthropometry Working Group,  
487 NHLBI Trans-Omics for Precision Medicine (TOPMed) Consortium, Cupples,  
488 L.A., Shadyab, A.H., McKnight, B., Shoemaker, B.M., Mitchell, B.D., Psaty,  
489 B.M., Kooperberg, C., et al.: Assessing the contribution of rare variants to com-  
490 plex trait heritability from whole-genome sequence data **54**, 263–273 (2022)  
491 <https://doi.org/10.1038/s41588-021-00997-7>
- 492 [5] Kosmicki, J.A., Churchhouse, C.L., Rivas, M.A., Neale, B.M.: Discovery of rare  
493 variants for complex phenotypes. *Hum. Genet.* **135**, 625–634 (2016) [https://doi.](https://doi.org/10.1007/s00439-016-1679-1)  
494 [org/10.1007/s00439-016-1679-1](https://doi.org/10.1007/s00439-016-1679-1)
- 495 [6] Fiziev, P.P., McRae, J., Ulirsch, J.C., Dron, J.S., Hamp, T., Yang, Y., Wain-  
496 schtein, P., Ni, Z., Schraiber, J.G., Gao, H., Cable, D., Field, Y., Aguet,  
497 F., Fasnacht, M., Metwally, A., Rogers, J., Marques-Bonet, T., Rehm, H.L.,  
498 O'Donnell-Luria, A., Khera, A.V., Farh, K.K.-H.: Rare penetrant mutations confer  
499 severe risk of common diseases. *Science* **380**, 1131 (2023) [https://doi.org/10.](https://doi.org/10.1126/science.abo1131)  
500 [1126/science.abo1131](https://doi.org/10.1126/science.abo1131)
- 501 [7] Hoogmartens, J., Cacace, R., Van Broeckhoven, C.: Insight into the genetic eti-  
502 ology of alzheimer's disease: A comprehensive review of the role of rare variants  
503 **13**, 12155 (2021) <https://doi.org/10.1002/dad2.12155>
- 504 [8] Li, Z., Li, X., Zhou, H., Gaynor, S.M., Selvaraj, M.S., Arapoglou, T., Quick, C.,  
505 Liu, Y., Chen, H., Sun, R., et al.: A framework for detecting noncoding rare-  
506 variant associations of large-scale whole-genome sequencing studies **19**, 1599–1611  
507 (2022) <https://doi.org/10.1038/s41592-022-01640-x>
- 508 [9] Lee, S., Abecasis, G.R., Boehnke, M., Lin, X.: Rare-variant association analysis:  
509 study designs and statistical tests. *Am. J. Hum. Genet.* **95**, 5–23 (2014) [https:](https://doi.org/10.1016/j.ajhg.2014.06.009)  
510 [//doi.org/10.1016/j.ajhg.2014.06.009](https://doi.org/10.1016/j.ajhg.2014.06.009)
- 511 [10] Li, X., Li, Z., Zhou, H., Gaynor, S.M., Liu, Y., Chen, H., Sun, R., Dey, R.,  
512 Arnett, D.K., Aslibekyan, S., et al.: Dynamic incorporation of multiple in silico  
513 functional annotations empowers rare variant association analysis of large whole-  
514 genome sequencing studies at scale **52**, 969–983 (2020) [https://doi.org/10.1038/](https://doi.org/10.1038/s41588-020-0676-4)  
515 [s41588-020-0676-4](https://doi.org/10.1038/s41588-020-0676-4)
- 516 [11] Clarke, B., Holtkamp, E., Öztürk, H., Mück, M., Wahlberg, M., Meyer, K.,  
517 Munzlinger, F., Brechtmann, F., Hölzlwimmer, F.R., Gagneur, J., Stegle, O.:

- 518 Integration of variant annotations using deep set networks boosts rare variant  
519 association genetics, 2023–0712548506 (2023) [https://doi.org/10.1101/2023.07.](https://doi.org/10.1101/2023.07.12.548506)  
520 [12.548506](https://doi.org/10.1101/2023.07.12.548506)
- 521 [12] Wagner, N., Çelik, M.H., Hölzlwimmer, F.R., Mertes, C., Prokisch, H., Yépez,  
522 V.A., Gagneur, J.: Aberrant splicing prediction across human tissues **55**, 861–870  
523 (2023) <https://doi.org/10.1038/s41588-023-01373-3>
- 524 [13] Zeng, T., Li, Y.I.: Predicting RNA splicing from DNA sequence using pangolin  
525 **23**, 103 (2022) <https://doi.org/10.1186/s13059-022-02664-4>
- 526 [14] Avsec, , Agarwal, V., Visentin, D., Ledsam, J.R., Grabska-Barwinska, A., Taylor,  
527 K.R., Assael, Y., Jumper, J., Kohli, P., Kelley, D.R.: Effective gene expression  
528 prediction from sequence by integrating long-range interactions. *Nat. Methods*  
529 **18**, 1196–1203 (2021) <https://doi.org/10.1038/s41592-021-01252-x>
- 530 [15] Pollard, K.S., Hubisz, M.J., Rosenbloom, K.R., Siepel, A.: Detection of non-  
531 neutral substitution rates on mammalian phylogenies **20**, 110–121 (2010) <https://doi.org/10.1101/gr.097857.109>  
532
- 533 [16] Ionita-Laza, I., McCallum, K., Xu, B., Buxbaum, J.D.: A spectral approach inte-  
534 grating functional genomic annotations for coding and noncoding variants **48**,  
535 214–220 (2016) <https://doi.org/10.1038/ng.3477>
- 536 [17] Skene, N.G., Grant, S.G.N.: Identification of vulnerable cell types in major  
537 brain disorders using single cell transcriptomes and expression weighted cell type  
538 enrichment **10**, 16 (2016) <https://doi.org/10.3389/fnins.2016.00016>
- 539 [18] Li, Z., Li, X., Liu, Y., Shen, J., Chen, H., Zhou, H., Morrison, A.C., Boerwinkle,  
540 E., Lin, X.: Dynamic scan procedure for detecting rare-variant association regions  
541 in whole-genome sequencing studies. *Am. J. Hum. Genet.* **104**, 802–814 (2019)  
542 <https://doi.org/10.1016/j.ajhg.2019.03.002>
- 543 [19] Sey, N.Y.A., Hu, B., Mah, W., Fauni, H., McAfee, J.C., Rajarajan, P., Bren-  
544 nand, K.J., Akbarian, S., Won, H.: A computational tool (H-MAGMA) for  
545 improved prediction of brain-disorder risk genes by incorporating brain chromatin  
546 interaction profiles. *Nat. Neurosci.* **23**, 583–593 (2020) [https://doi.org/10.1038/](https://doi.org/10.1038/s41593-020-0603-0)  
547 [s41593-020-0603-0](https://doi.org/10.1038/s41593-020-0603-0)
- 548 [20] Ma, S., Dagleish, J., Lee, J., Wang, C., Liu, L., Gill, R., Buxbaum, J.D., Chung,  
549 W.K., Aschard, H., Silverman, E.K., Cho, M.H., He, Z., Ionita-Laza, I.: Powerful  
550 gene-based testing by integrating long-range chromatin interactions and knockoff  
551 genotypes **118** (2021) <https://doi.org/10.1073/pnas.2105191118>
- 552 [21] Zhang, S., Moll, T., Rubin-Sigler, J., Tu, S., Li, S., Yuan, E., Liu, M., Butt, A.,  
553 Harvey, C., et al.: Deep learning modeling of rare noncoding genetic variants in  
554 human motor neurons defines *CCDC146* as a therapeutic target for ALS (2024)

- 555 <https://doi.org/10.1101/2024.03.30.24305115>
- 556 [22] Fulco, C.P., Nasser, J., Jones, T.R., Munson, G., Bergman, D.T., Subrama-  
557 nian, V., Grossman, S.R., Anyoha, R., Doughty, B.R., Patwardhan, *et al.*:  
558 Activity-by-contact model of enhancer-promoter regulation from thousands of  
559 CRISPR perturbations. *Nat. Genet.* **51**, 1664–1669 (2019) <https://doi.org/10.1038/s41588-019-0538-0>  
560
- 561 [23] Hoffman, M.D., Blei, D.M., Wang, C., Paisley, J.: Stochastic variational inference.  
562 *J. Mach. Learn. Res.* (2013)
- 563 [24] Jaganathan, K., Kyriazopoulou Panagiotopoulou, S., McRae, J.F., Darbandi,  
564 S.F., Knowles, D., Li, Y.I., Kosmicki, J.A., Arbelaez, J., Cui, W., Schwartz,  
565 G.B., Chow, E.D., Kanterakis, E., Gao, H., Kia, A., Batzoglou, S., Sanders, S.J.,  
566 Farh, K.K.-H.: Predicting splicing from primary sequence with deep learning **176**,  
567 535–54824 (2019) <https://doi.org/10.1016/j.cell.2018.12.015>
- 568 [25] Madsen, B.E., Browning, S.R.: A groupwise association test for rare mutations  
569 using a weighted sum statistic. *PLoS Genet.* **5**, 1000384 (2009) <https://doi.org/10.1371/journal.pgen.1000384>  
570
- 571 [26] Morgenthaler, S., Thilly, W.G.: A strategy to discover genes that carry multi-  
572 allelic or mono-allelic risk for common diseases: a cohort allelic sums test (CAST)  
573 **615**, 28–56 (2007) <https://doi.org/10.1016/j.mrfmmm.2006.09.003>
- 574 [27] He, Z., Xu, B., Lee, S., Ionita-Laza, I.: Unified sequence-based association tests  
575 allowing for multiple functional annotations and meta-analysis of noncoding vari-  
576 ation in metabochip data **101**, 340–352 (2017) <https://doi.org/10.1016/j.ajhg.2017.07.011>  
577
- 578 [28] Nott, A., Holtman, I.R., Coufal, N.G., Schlachetzki, J.C.M., Yu, M., Hu, R., Han,  
579 C.Z., Pena, M., Xiao, J., Wu, *et al.*: Brain cell type-specific enhancer-promoter  
580 interactome maps and disease-risk association. *Science* **366**, 1134–1139 (2019)  
581 <https://doi.org/10.1126/science.aay0793>
- 582 [29] Dempster, E.R., Lerner, I.M.: Heritability of threshold characters **35**, 212–236  
583 (1950) <https://doi.org/10.1093/genetics/35.2.212>
- 584 [30] Zhou, X., Chen, Y., Ip, F.C.F., Jiang, Y., Cao, H., Lv, G., Zhong, H., Chen, J.,  
585 Ye, T., Chen, *et al.*: Deep learning-based polygenic risk analysis for alzheimer’s  
586 disease prediction. *Commun. Med. (Lond.)* **3**, 49 (2023) <https://doi.org/10.1038/s43856-023-00269-x>  
587
- 588 [31] Blei, D.M., Kucukelbir, A., McAuliffe, J.D.: Variational inference: A review for  
589 statisticians. *J. Am. Stat. Assoc.* **112**, 859–877 (2017) <https://doi.org/10.1080/01621459.2017.1285773>  
590

- 591 [32] Gopinath, A., Collins, A., Khoshbouei, H., Streit, W.J.: Microglia and other  
592 myeloid cells in central nervous system health and disease. *J. Pharmacol. Exp.*  
593 *Ther.* **375**, 154–160 (2020) <https://doi.org/10.1124/jpet.120.265058>
- 594 [33] Gratuze, M., Leyns, C.E.G., Holtzman, D.M.: New insights into the role of  
595 TREM2 in alzheimer’s disease. *Mol. Neurodegener.* **13**, 66 (2018) [https://doi.](https://doi.org/10.1186/s13024-018-0298-9)  
596 [org/10.1186/s13024-018-0298-9](https://doi.org/10.1186/s13024-018-0298-9)
- 597 [34] Han, Y., Chen, K., Yu, H., Cui, C., Li, H., Hu, Y., Zhang, B., Li, G.: Maf1 loss  
598 regulates spinogenesis and attenuates cognitive impairment in alzheimer’s disease.  
599 *Brain* **147**, 2128–2143 (2024) <https://doi.org/10.1093/brain/awae015>
- 600 [35] Bottero, V., Potashkin, J.A.: Meta-analysis of gene expression changes in  
601 the blood of patients with mild cognitive impairment and alzheimer’s disease  
602 dementia. *Int. J. Mol. Sci.* **20**, 5403 (2019) <https://doi.org/10.3390/ijms20215403>
- 603 [36] Izadi, F., Soheilifar, M.H.: Exploring potential biomarkers underlying patho-  
604 genesis of alzheimer’s disease by differential co-expression analysis **10**, 233–241  
605 (2018)
- 606 [37] Xiao, B., Kuang, Z., Zhang, W., Hang, J., Chen, L., Lei, T., He, Y., Deng,  
607 C., Li, W., Lu, J., Qu, J., et al.: Glutamate ionotropic receptor kainate type  
608 subunit 3 (GRIK3) promotes epithelial-mesenchymal transition in breast can-  
609 cer cells by regulating SPDEF/CDH1 signaling **58**, 1314–1323 (2019) [https:](https://doi.org/10.1002/mc.23014)  
610 [//doi.org/10.1002/mc.23014](https://doi.org/10.1002/mc.23014)
- 611 [38] Spead, O., Zaepfel, B.L., Rothstein, J.D.: Nuclear pore dysfunction in neurode-  
612 generation. *Neurotherapeutics* **19**, 1050–1060 (2022) [https://doi.org/10.1007/](https://doi.org/10.1007/s13311-022-01293-w)  
613 [s13311-022-01293-w](https://doi.org/10.1007/s13311-022-01293-w)
- 614 [39] Bukvic, N., De Rinaldis, M., Chetta, M., Trabacca, A., Bassi, M.T., Marsano,  
615 R.M., Holoubkova, L., Riveccio, M., Oro, M., Resta, *et al.*: De novo pathogenic  
616 variant in FBRSL1, non OMIM gene paralogue AUTS2, causes a novel recogniz-  
617 able syndromic manifestation with intellectual disability; an additional patient  
618 and review of the literature. *Genes (Basel)* **15**, 826 (2024) [https://doi.org/10.](https://doi.org/10.3390/genes15070826)  
619 [3390/genes15070826](https://doi.org/10.3390/genes15070826)
- 620 [40] Mishra, R., Li, B.: The application of artificial intelligence in the genetic study of  
621 alzheimer’s disease. *Aging Dis.* **11**, 1567–1584 (2020) [https://doi.org/10.14336/](https://doi.org/10.14336/AD.2020.0312)  
622 [AD.2020.0312](https://doi.org/10.14336/AD.2020.0312)
- 623 [41] Lutz, M.W., Sprague, D., Barrera, J., Chiba-Falek, O.: Shared genetic etiology  
624 underlying alzheimer’s disease and major depressive disorder. *Transl. Psychiatry*  
625 **10**, 88 (2020) <https://doi.org/10.1038/s41398-020-0769-y>
- 626 [42] GTEx Consortium: The GTEx consortium atlas of genetic regulatory effects  
627 across human tissues **369**, 1318–1330 (2020) <https://doi.org/10.1126/science>.

628 aaz1776

- 629 [43] Humphrey, J., Brophy, E., Kosoy, R., Zeng, B., Coccia, E., Mattei, D., Ravi,  
630 A., Efthymiou, A.G., Navarro, E., Muller, B.Z., et al.: Long-read RNA-seq atlas  
631 of novel microglia isoforms elucidates disease-associated genetic regulation of  
632 splicing. *medRxiv*, 2023–120123299073 (2023) <https://doi.org/10.1101/2023.12.01.23299073>  
633
- 634 [44] Bellenguez, C., Küçükali, F., Jansen, I.E., Kleindam, L., Moreno-Grau, S.,  
635 Amin, N., Naj, A.C., Campos-Martin, R., Grenier-Boley, B., Andrade, *et al.*: New  
636 insights into the genetic etiology of alzheimer’s disease and related dementias.  
637 *Nat. Genet.* **54**, 412–436 (2022) <https://doi.org/10.1038/s41588-022-01024-z>
- 638 [45] Lambert, J.C., Ibrahim-Verbaas, C.A., Harold, D., Naj, A.C., Sims, R., Bel-  
639 lenguez, C., DeStafano, A.L., Bis, J.C., Beecham, G.W., Grenier-Boley, B., Russo,  
640 *et al.*: Meta-analysis of 74,046 individuals identifies 11 new susceptibility loci for  
641 alzheimer’s disease. *Nat. Genet.* **45**, 1452–1458 (2013) [https://doi.org/10.1038/](https://doi.org/10.1038/ng.2802)  
642 [ng.2802](https://doi.org/10.1038/ng.2802)
- 643 [46] Nalls, M.A., Blauwendraat, C., Vallerga, C.L., Heilbron, K., Bandres-Ciga, S.,  
644 Chang, D., Tan, M., Kia, D.A., Noyce, A.J., *et al.*: Identification of novel risk  
645 loci, causal insights, and heritable risk for parkinson’s disease: a meta-analysis  
646 of genome-wide association studies. *Lancet Neurol.* **18**, 1091–1102 (2019) [https://doi.org/10.1016/S1474-4422\(19\)30320-5](https://doi.org/10.1016/S1474-4422(19)30320-5)  
647
- 648 [47] Desikan, R.S., Schork, A.J., Wang, Y., Witoelar, A., Sharma, M., McEvoy, L.K.,  
649 Holland, D., Brewer, J.B., Chen, C.-H., Thompson, W.K., *et al.*: Genetic overlap  
650 between alzheimer’s disease and parkinson’s disease at the MAPT locus. *Mol.*  
651 *Psychiatry* **20**, 1588–1595 (2015) <https://doi.org/10.1038/mp.2015.6>
- 652 [48] Kuksa, P.P., Liu, C.-L., Fu, W., Qu, L., Zhao, Y., Katanic, Z., Clark, K., Kuzma,  
653 A.B., Ho, P.-C., *et al.*: Alzheimer’s disease variant portal: A catalog of genetic  
654 findings for alzheimer’s disease. *J. Alzheimers. Dis.* **86**, 461–477 (2022) <https://doi.org/10.3233/JAD-215055>  
655
- 656 [49] Wu, M.C., Lee, S., Cai, T., Li, Y., Boehnke, M., Lin, X.: Rare-variant association  
657 testing for sequencing data with the sequence kernel association test. *Am. J. Hum.*  
658 *Genet.* **89**, 82–93 (2011) <https://doi.org/10.1016/j.ajhg.2011.05.029>
- 659 [50] Lee, S., Emond, M.J., Bamshad, M.J., Barnes, K.C., Rieder, M.J., Nickerson,  
660 D.A., NHLBI GO Exome Sequencing Project—ESP Lung Project Team, Chris-  
661 tiani, D.C., Wurfel, M.M., Lin, X.: Optimal unified approach for rare-variant  
662 association testing with application to small-sample case-control whole-exome  
663 sequencing studies **91**, 224–237 (2012) [https://doi.org/10.1016/j.ajhg.2012.06.](https://doi.org/10.1016/j.ajhg.2012.06.007)  
664 [007](https://doi.org/10.1016/j.ajhg.2012.06.007)
- 665 [51] Liu, Y., Chen, S., Li, Z., Morrison, A.C., Boerwinkle, E., Lin, X.: ACAT: A fast

- 666 and powerful p value combination method for rare-variant analysis in sequencing  
667 studies **104**, 410–421 (2019) <https://doi.org/10.1016/j.ajhg.2019.01.002>
- 668 [52] Beecham, G.W., Bis, J.C., Martin, E.R., Choi, S.-H., DeStefano, A.L., Duijn,  
669 C.M., Fornage, M., Gabriel, S.B., Koboldt, D.C., Larson, D.E., *et al.*: The  
670 alzheimer’s disease sequencing project: Study design and sample selection. *Neurol*  
671 *Genet* **3**, 194 (2017) <https://doi.org/10.1212/NXG.000000000000194>
- 672 [53] Manichaikul, A., Mychaleckyj, J.C., Rich, S.S., Daly, K., Sale, M., Chen, W.-M.:  
673 Robust relationship inference in genome-wide association studies **26**, 2867–2873  
674 <https://doi.org/10.1093/bioinformatics/btq559>
- 675 [54] Purcell, S., Neale, B., Todd-Brown, K., Thomas, L., Ferreira, M.A.R., Bender, D.,  
676 Maller, J., Sklar, P., Bakker, P.I.W., Daly, M.J., Sham, P.C.: PLINK: a tool set  
677 for whole-genome association and population-based linkage analyses **81**, 559–575  
678 (2007) <https://doi.org/10.1086/519795>
- 679 [55] Siva, N.: 1000 genomes project **26**, 256 (2008) <https://doi.org/10.1038/nbt0308-256b>
- 681 [56] Ge, T., Chen, C.-Y., Ni, Y., Feng, Y.-C.A., Smoller, J.W.: Polygenic prediction  
682 via bayesian regression and continuous shrinkage priors. *Nat. Commun.* **10**, 1776  
683 (2019) <https://doi.org/10.1038/s41467-019-09718-5>
- 684 [57] Eli, B., Jonathan, P.C., Martin, J., Fritz, O., Neeraj, P., Theofanis, K., Rohit, S.,  
685 Paul, S., Paul, H., Noah, D.G.: Pyro: Deep universal probabilistic programming.  
686 arXiv [cs.LG] (2018) <https://doi.org/10.48550/arXiv.1810.09538>
- 687 [58] Liu, X., White, S., Peng, B., Johnson, A.D., Brody, J.A., Li, A.H., Huang, Z.,  
688 Carroll, A., Wei, P., Gibbs, R., Klein, R.J., Boerwinkle, E.: WGSAs: an annotation  
689 pipeline for human genome sequencing studies **53**, 111–112 (2016) <https://doi.org/10.1136/jmedgenet-2015-103423>
- 691 [59] Bernstein, B.E., Stamatoyannopoulos, J.A., Costello, J.F., Ren, B., Milosavljevic,  
692 A., Meissner, A., Kellis, M., Marra, M.A., Beaudet, A.L., Ecker, J.R., *et al.*:  
693 The NIH roadmap epigenomics mapping consortium **28**, 1045–1048 (2010) <https://doi.org/10.1038/nbt1010-1045>
- 695 [60] Chen, S., Francioli, L.C., Goodrich, J.K., Collins, R.L., Kanai, M., Wang, Q.,  
696 Alföldi, J., Watts, N.A., Vittal, C., Gauthier, L.D., *et al.*: A genomic mutational  
697 constraint map using variation in 76,156 human genomes. *Nature* **625**, 92–100  
698 (2024) <https://doi.org/10.1038/s41586-023-06045-0>
- 699 [61] Karczewski, K.J., Francioli, L.C., Tiao, G., Cummings, B.B., Alföldi, J., Wang,  
700 Q., Collins, R.L., Laricchia, K.M., Ganna, A., Birnbaum, D.P., Gauthier, *et al.*:  
701 The mutational constraint spectrum quantified from variation in 141,456 humans.  
702 *Nature* **581**, 434–443 (2020) <https://doi.org/10.1038/s41586-020-2308-7>

- 703 [62] Weissbrod, O., Hormozdiari, F., Benner, C., Cui, R., Ulirsch, J., Gazal, S.,  
704 Schoech, A.P., Geijn, B., Reshef, Y., Márquez-Luna, C., O'Connor, L., Pirinen,  
705 M., Finucane, H.K., Price, A.L.: Functionally informed fine-mapping and poly-  
706 genic localization of complex trait heritability. *Nat Genet* **52**, 1355–1363 (2020)  
707 <https://doi.org/10.1038/s41588-020-00735-5>
- 708 [63] Szustakowski, J.D., Balasubramanian, S., Kvikstad, E., Khalid, S., Bronson,  
709 P.G., Sasson, A., Wong, E., Liu, D., Wade Davis, J., Haefliger, *et al.*: Advanc-  
710 ing human genetics research and drug discovery through exome sequencing  
711 of the UK biobank. *Nat Genet* **53**, 942–948 (2021) <https://doi.org/10.1038/s41588-021-00885-0>
- 713 [64] Zhou, W., Zhao, Z., Nielsen, J.B., Fritsche, L.G., LeFaive, J., Gagliano Taliun,  
714 S.A., Bi, W., Gabrielsen, M.E., Daly, M.J., Neale, B.M., et al.: Scalable general-  
715 ized linear mixed model for region-based association tests in large biobanks and  
716 cohorts **52**, 634–639 (2020) <https://doi.org/10.1038/s41588-020-0621-6>
- 717 [65] Londhe, S., Lindner, J., Chen, Z., Holtkamp, E., Hölzlwimmer, F.R., Casale, F.P.,  
718 Brechtmann, F., Gagneur, J.: Functional gene embeddings improve rare variant  
719 polygenic risk scores. *bioRxiv*, 2024–0722604535 (2024) <https://doi.org/10.1101/2024.07.22.604535>  
720