

1 **Title:**

2 Biological machine learning combined with bacterial population genomics reveals common
3 and rare allelic variants of genes to cause disease

4

5 **Authors:**

6 DJ Darwin R. Bandoy^{1,2} and Bart C. Weimer^{1*}

7

8 **Affiliations:**

9 ¹University of California Davis, School of Veterinary Medicine, 100 K Pathogen Genome Project,
10 Davis, CA 95616, USA; ²University of the Philippines Los Baños, College of Veterinary
11 Medicine, Department of Veterinary Paraclinical Sciences, Laguna 4031, Philippines

12

13 *corresponding author: bcweimer@ucdavis.edu; 530-760-9550

14

15

16

17 **Key words:**

18 Infectious disease, XGboost, *Campylobacter*, abortion, protein modeling

19 **Abstract**

20 Highly dimensional data generated from bacterial whole genome sequencing is providing
21 unprecedented scale of information that requires appropriate statistical frameworks of analysis
22 to infer biological function from bacterial genomic populations. Application of genome wide
23 association study (GWAS) methods is an emerging approach with bacterial population
24 genomics that yields a list of genes associated with a phenotype with an undefined importance
25 among the candidates in the list. Here, we validate the combination of GWAS, machine
26 learning, and pathogenic bacterial population genomics as a novel scheme to identify SNPs and
27 rank allelic variants to determine associations for accurate estimation of disease phenotype.
28 This approach parsed a dataset of 1.2 million SNPs that resulted in a ranked importance of
29 associated alleles of *Campylobacter jejuni porA* using multiple spatial locations over a 30-year
30 period. We validated this approach using previously proven laboratory experimental alleles from
31 an *in vivo* guinea pig abortion model. This approach, termed BioML, defined intestinal and
32 extraintestinal groups that have differential allelic variants that cause abortion. Divergent
33 variants containing indels that defeated gene callers were rescued using biological context and
34 knowledge that resulted in defining rare and divergent variants that were maintained in the
35 population over two continents and 30 years. This study defines the capability of machine
36 learning coupled to GWAS and population genomics to simultaneously identify and rank alleles
37 to define their role in abortion, and more broadly infectious disease.

38 39 **Main**

40 Comparative microbial genomics has emerged from pangenome comparisons that are
41 exclusively tied to reference genomes that define the perspective of change to a core and
42 flexible genome perspective lacking a firm confirmation of which genes are linked to disease¹.
43 An alternative approach to this perspective is use of genome wide association (GWAS) methods
44 that are common in mammalian genomics in an effort to refine the estimates of specific genes of

45 interest. A limitation of GWAS is that it sequentially examines single loci that prevents
46 simultaneous analysis of different allelic variants that can be interacting at different levels and
47 population distribution between strain differentiation². This is a severe limitation in bacterial
48 genomics, especially as population genomics is now possible in bacteria at a scale that allows
49 examination of non-linear evolutionary rates of each gene and all of the alleles found in very
50 large populations that create big data analytical problems. A compounding limitation is the lack
51 of appropriate statistical models that underpin this approach in bacteria since it is unknown
52 when the populations are normally distributed or evolving in a non-linear progression. As with all
53 large data sets, multiple comparisons require Bonferroni correction to adjust the p -value based
54 on a new scale as compared to gene expression but it is on a scale that is beyond that
55 contemplated for gene expression variation (Table 1)³. Further, the assumption that each gene
56 or allele is independent is conceptually flawed; and hence, alternative analyses that are
57 biological and statistically compatible needs to be defined.

58 Coupling GWAS, population microbial genomics, and machine learning is poised to be a
59 robust alternative to classical GWAS or pangenome comparison to simultaneously discover
60 changes in microbial genomes, and genes, that span the scale of genome plasticity to alleles of
61 a single gene. Moreover, this combination (coined BioML) will produce a statistically
62 underpinned comparative ranking of the most important factors that are not obvious from GWAS
63 alone. These advantages combined with downstream inspection of the prioritized rank further
64 powers discovery to bring biologically insightful observations and solutions, especially when
65 large genome populations are used in the analysis, from very divergent populations of alleles
66 that are missed when gene calling is too divergent.

67 An analytical strength for use of machine learning in microbiology is the ability to define
68 functional relationship from population scale genomes or genes without a priori definition of the
69 underlying mechanism of change or specific phenotype limitations⁴. This distinctive advantage
70 makes machine learning superior to classical statistical tests for prokaryotic systems that are

71 highly variable, particularly bacteria wherein explanatory variables are not linearly correlated,
72 features are dependent due to genome linkage, varying evolutionary rates of between genes,
73 and assumptions of normal distribution are violated in part due to varying selection conditions^{2,5}.
74 These biological conditions and parameters are incompatible with the assumptions of linear or
75 correlative statistics, which is compounded with data reduction methods that provide a very
76 small snapshot of the genome variation that yield associations that have low predictive value.

77 In this study, we verified the concept of coupling GWAS with machine learning and
78 population bacterial genomics (Figure 1) in a use case to test the hypothesis that a specific
79 gene (*porA*) is linked to extraintestinal location and further is causative in abortion⁶⁻⁸ in a ranked
80 order that is biologically meaningful. This was done using a wet lab validated data set containing
81 100 genomes⁶⁻⁸ using extreme gradient boosting (XGboost), which was used in biological
82 applications previously⁹. XGboost can identify genetic variants in human GWAS as
83 demonstrated in a Finnish study that integrated complex nonlinear interactions of SNPs¹⁰. The
84 ability to interrogate the predictive features enables whiteboxing the parameters, which is
85 emerging as a tool for deriving mechanistic function in biology¹¹. XGboost implements adaptive
86 optimization within the functional space by iteration of the weak learners into strong learners
87 represented by decision trees where each new decision tree is generated by factoring the
88 residual generated from the difference from observed to the predicted feature (Figure 2;
89 Supplemental Table 1).

90 This study used a previously validated wet lab data set with a tetracycline resistant strain of
91 *Campylobacter jejuni* causing abortion in sheep⁶⁻⁸. Their studies used a pairwise genome
92 comparison to identify 8,000 SNP difference between a reference genome and abortive strain
93 and utilized transformed genomes to identify specific allelic variants driving abortion. We utilized
94 those 85 genomes that span 30 years and multiple locations as a reference set of cases and
95 108 control genomes of intestinal, diarrheal isolates. This approach allows exploration of
96 bacterial population genomic space by linking different phenotypes to the genome variation

97 among the isolates (Figure 1). Biological feature engineering of this collection of genomes
98 identified 1.2 million SNPs, which is not tractable using in vivo infection studies to determine the
99 roll of all SNPs. To examine this scale problem, we hypothesized that genomic changes evolved
100 in gastrointestinal *C. jejuni* resulting in an abortive phenotype; hence, moving from the intestine
101 to other tissues – in this case the placenta resulting in abortion. Applying our approach (BioML)
102 to a population of gastrointestinal, diarrheal *C. jejuni* versus extraintestinal, abortive phenotypes
103 produced a prioritized allelic difference in a ranked order of importance to the phenotype (i.e.
104 abortion) (Supplementary Table 1).

105 BioML identified 14 *porA* loci as the most important alleles ranging from 89 to 59 relative
106 importance out of the 1.2 million SNPs (Supplemental Table 1). These ranked loci were
107 compared by body location (Figure 3), which further clarified the location of these SNPs in a
108 Tetris plot that simultaneously presented the ranked associated allelic variants within the
109 phenotype of interest as detected with BioML as well as the non-associated alleles. By
110 presenting the cases and control simultaneously within the y-axis, capturing insight is easily
111 observed and areas for further investigation can be prioritized with visual inspection combined
112 with biological knowledge. An added feature of the Tetris plot, which is lacking in Manhattan
113 plots, is the ability to detect rare variants that are not captured by gene calling, machine learning
114 alone, or classical statistical testing. Regions within cases expressing different allelic patterns
115 were further explored for each genome and implications in biological features important in the
116 disease. Additionally, protein structures were modeled to examine the changes in protein
117 configuration initially yielded three distinct groups (Figure 3). These alleles were directly
118 compared to those validated *in vivo* and found to be linked to specific protein loops within alleles
119 verified previously⁶⁻⁸ – BioML found each of those to be biologically important for abortion.

120 Further we located each of the top ranked alleles loops 1, 3, 4, 7 as enriched selection loci,
121 again verifying previous wet lab observations⁶⁻⁸. Tetris plot derived variants that were not 100%
122 identical with >75% protein homology, were designated as nonprototypical variants because the

123 sequence variation was high enough to change protein structures. In a limited set of alleles, the
124 *porA* allele was so divergent that they were not variant called but were recovered with manual
125 curation for this study. Recovery of these genes that were not initially identified created a third
126 group of rare variant alleles that also caused abortion (Figure 4; protein homology <75 %). All of
127 the variants were mapped to the whole genome phylogeny diversity as well rare variants that
128 were not variant called by the reference genome (Figure 4). Prototypical allelic variants
129 clustered in the largest genomic group of abortive isolates, as did some of the nonprototypical
130 *porA* variants. However, there was significant genome variation and contained the two groups
131 that caused abortion. Rare *porA* variants were distributed within different genomic groups as
132 well as over a 15-year span between North America and the UK. The extensive allelic versions
133 of *porA*, as well as the different genotypes, suggests that a genome surveillance system based
134 on SNPs would be unsuccessful to link these genomes to a disease. In combination, these
135 observations indicate that BioML produced a ranked list of biologically important alleles that
136 were validated with those that were previously shown to be causal in abortion for the exact SNP
137 and the protein loop location. Together, these observations verified that BioML was capable of
138 accurately identifying the exact SNPs in *porA* that cause abortion.

139 Since each BioML allele was validated for accuracy to wet lab results correctly, we further
140 examined the protein changes from the ranked alleles (Figure 5). The first six top ranked alleles
141 changed the amino acids for each *porA* sequence, but each protein sequenced varied across all
142 PorA models. However, lysine₁₈₉ was conserved across the extraintestinal variants and Asn was
143 found in the intestinal alleles. Lysine mutation changes are the most impactful in membrane
144 pore structure and are one of the tenets of membrane topology as positive inside rule^{12,13}.
145 Positive inside rule describes the observation across membrane pores that positively charged
146 amino acids are found within the cytoplasm and negatively charged amino acids are in the
147 extracellular domain. Membrane topology can radically change from being oriented inside the
148 membrane (exposed to the periplasm in this case) to outside the membrane with a single lysine

149 mutation. Within the adjacent protein structure, lysine snorkeling effectively minimizes the
150 nonpolar chain component by burying in the hydrophobic domain and at the same time expose
151 the polar component to the aqueous domain is another single amino acid change which alters
152 the topology of the membrane domain¹⁴. Bacterial membrane pore flipping could be a potential
153 mechanism to avoid recognition by the immune system and enhancement of ion transport.
154 While the counterpart position is buried in a deeper position due to insertional mutation in rare
155 variants, the inserted amino acids contain lysine at new position 197. Additionally, insertions in
156 the rare variants reduce the homology to < 75 % lead to more extensive protein structural
157 changes that change the PorA arrangement in the membrane while still able to cause abortion.
158 This situation is troublesome for traditional approaches but BioML effectively identified this
159 situation.

160 This study utilized a combination of GWAS, population bacterial genomics, and machine
161 learning to identify and rank allelic variants that correspond to biologically validated alleles of
162 *porA* to cause abortion. BioML results were further supported by the longitudinal and spatial
163 conservation of *porA* coupled to protein substitutions that led to biologically relevant changes in
164 the structure to change activity. A Tetris plot visualization provided an avenue to discover
165 divergent and rare variants that provided further insight with protein modelling that uncovered
166 protein substitutions resulting in localization changes that affect activity and isolation localization
167 in the host. Together these results demonstrate and validate a novel method, termed BioML, to
168 discover biological mechanisms using population bacterial genomics. This approach provides
169 an avenue to leverage the massive amount of bacterial genomic sequences to uncover new
170 mechanisms of disease with potential to provide therapeutic approaches.

171

172 **Acknowledgement:**

173 DDB is grateful for the funding provided by the Philippine California Advanced Research

174 Institute and University of the Philippines Enhanced Creative Work and Research Grant to fund

175 his PhD program.

176

177 **Figure legends:**

178 **Figure 1.** Biological feature engineering of genomic data for machine learning analysis. A

179 critical step in feature engineering is selection of the appropriate comparison groups to enable

180 classification of alleles that are related to the specific phenotype of interest (i.e. intestinal

181 (controls; diarrheal; n=108) and extraintestinal (cases; abortive; n=85) (Step 1). Population-wide

182 allelic variants (red dot = intestinal, green dot = extraintestinal) that result from variant calling

183 (Step 2) and are used as the input features for machine learning analysis (Step 3). The

184 predicted model generated from the machine learning analysis is inspected for the most

185 predictive features using biological context, input, and protein modelling (Step 4) that represents

186 a nonsynonymous mutation from the genomic the population of allelic variants (n=193).

187

188 **Figure 2.** The conceptual framework diagram depicting machine learning in bacterial genome

189 wide association using extreme gradient boosting (XGboost). Boosting is a technique of

190 combining a set of weak classifiers or decision trees to increase prediction accuracy. Red dots

191 represent an allelic variant, each grey bar represents a unique allele. Individual decision trees

192 (1, 2, 3) fail to fully capture the allelic variants associated with the phenotype (e.g. extraintestinal

193 abortion), but by combining the trees together results in a process called as boosting increases

194 the discriminative power.

195

196 **Figure 3.** Comparative plot of SNP loci along the *proA* gene in all genomes. We termed this a

197 Tetris plot as an alternative visualization of genome wide association hits because they are

198 ranked and display only the loci that vary to produce a nonsynonymous mutation. The y-axis

199 contains individual genomes from the cases and the controls, while the x-axis contains the

200 GWAS SNP loci (green), the non-disease associated SNPs (red), open space (white) are loci

201 that are identical in the gene sequence. Temporal and geographic metadata on the right side of

202 the Tetris plot provides context for mutational enrichment over 30 years and multiple distant

203 locations in North America and the UK. The enriched SNP variation produced different protein

204 structures (far right in blue) as the corresponding protein model by location within the animal by
205 SNP. Protein structural features corresponding to the ranked GWAS variants are annotated on
206 top and below the plot are the nucleotide coordinates. Rare variants (homology <75%) was not
207 included by the variant caller in this visualization but manual inspection provided a method to
208 find these variants.

209
210 **Figure 4.** Whole genome distance matrix using minhash depicting an all against all comparison
211 of genome diversity for all isolates used in this study overlaid with the *porA* variant associated
212 with body location and disease phenotype. Genotypes and *porA* variants are connected in this
213 depiction to examine the association between intestinal/diarrheal location (yellow dot boxes),
214 prototypical extraintestinal/abortive (red dot boxes), non-prototypical *porA* variants in
215 extraintestinal/abortive (maroon lines), and rare *porA* variants in extraintestinal/abortive (grey
216 dashed lines) were co-located to their respective genomes in the genotype map. For the non-
217 prototypical variants, the year and location of isolation was included to depict the variation over
218 time and space in the maintenance of a minority population of *porA* variants of extraintestinal
219 abortive *Campylobacter jejuni*. The diagram to the right depicts the process used for this
220 analysis.

221
222 **Figure 5.** Protein models of the four groups of *porA* allelic variants that change the protein
223 model structure relative to the isolate location in the host and the disease outcome. The amino
224 acids corresponding with the BioML top ranked alleles are labelled in the common variant of
225 *porA*, while the rest show the substituted amino acid in their respective position.

226
227
228 **List of Tables**

229 Table 1. Exemplar comparison of statistical metrics of GWAS versus machine learning metrics.

230 Allelic variant association with phenotype using XGboost. An allele can be very large, ~8,000 for
231 *porA* for a pairwise comparison. Using a population of this gene from 200 genomes created a

232 population variation of 1.2 million variants that can be ranked with an estimation of importance

233 to association with the disease phenotype, abortion in this case.

234

235 **List of Supplemental Tables and Figures**

236 Supplemental Table 1. Ranked allelic variants using BioML

237 Supplemental Table 2. Metadata for extraintestinal *Campylobacter jejuni*

238 Supplemental Table 3. Metadata for intestinal *Campylobacter jejuni*

239 Supplemental Table 4. Confusion matrix and derived model metrics for the XGboost model with

240 extraintestinal *Campylobacter jejuni*. TP= True positive, FN= False Negative, FP= False

241 Positive, FN= False Negative.

242

243 **Methods**

244 **Biological feature engineering**

245 Biological feature engineering entails selection of pertinent controls and cases for BioML

246 analysis. The genomes between gastrointestinal and extraintestinal abortive isolates. *C. jejuni*

247 controls were downloaded from Patric 3.5.28 (<https://www.patricbrc.org/>), June 1, 2019

248 (Supplemental Table 2). Abortive extraintestinal genomes of *C. jejuni* were obtained from the

249 Sequence Read Archive (SRA; Supplemental Table 3)⁸. Fastq files were assembled using

250 Shovill 1.0.4 (<https://github.com/tseemann/shovill>). Assembled files were annotated with Prokka

251 (version 1.13.3)¹⁵. Variant calling was done with the reference sequence *C. jejuni* NTC11168

252 with Snippy 4.3.5 (<https://github.com/tseemann/snippy>) as previously described¹⁶.

253

254 **Gradient tree boosting as GWAS framework**

255 GWAS variants generated from the biological feature engineering step were used as input for

256 XGboost. The source code for implementing gradient tree boosting is available at

257 <https://xgboost.readthedocs.io/>. Confusion matrix were generated and used to assess the

258 performance of the model (Supplemental Table 4). The relative importance of the predictive

259 model was used as the GWAS hits.

260

261 **Tetris plot**

262 Classical GWAS hits are displayed as the negative logarithm of the p-value in Manhattan plots,
263 hence we formulated a novel visualization of the ranked alleles generated by the machine
264 learning model to highlight the difference between approaches - we call this GWAS hit
265 visualization a Tetris plot. We color coded the relative importance values of the associated
266 alleles derived from the XGboost (green being associated and red being non-associated). The
267 source genome is plotted on the y-axis and genomic coordinates on the x-axis overlaid with
268 GWAS hits presence or absence matrix.

269

270 **Population wide whole genome phylogeny**

271 The genome distance metric was calculated using genome wide k-mer signatures to generate
272 the population-wide phylogeny with a k-mer size of 31 scaled to 1000 with Sourmash¹⁷. The
273 resulting genome wide k-mer distance was visualized as an all-against-all heatmap¹⁷.

274

275 **Protein Modelling**

276 Assembled genomes were annotated using Prokka (V1.13.3) and PorA protein sequences were
277 extracted for protein modelling using Swiss Model^{18,19}. The most homologous protein was used
278 as template for protein modelling. Illustrate (<https://ccsb.scripps.edu/illustrate/>) was used to
279 generate the protein visualization of the predictive alleles. Ranked BioML alleles identified by
280 visual inspection of the Tetris plot, via the ranked variable importance were used to inspect the
281 protein structures.

282

283

284 Table 1.

	GWAS statistical metrics		Machine Learning coupled to GWAS metrics	
Allele	GWAS p-value	Bonferonni Corrected p-value	Candidate Ranking	Feature Importance
X_1	0.001	8.3×10^{-10}	1	80
X_2	0.001	8.3×10^{-10}	2	75
X_3	0.001	8.3×10^{-10}	3	70
X_n	0.001	8.3×10^{-10}	Rank _n	Importance _n

285
286
287

288
289
290
291
292
293
294
295
296
297
298
299
300
301
302
303
304
305
306
307
308
309
310
311
312
313
314
315
316
317
318
319
320
321
322
323
324
325
326
327
328
329
330
331

References

- 1 Page, A. J. *et al.* Roary: rapid large-scale prokaryote pan genome analysis. *Bioinformatics* **31**, 3691-3693, doi:10.1093/bioinformatics/btv421 (2015).
- 2 Power, R. A., Parkhill, J. & de Oliveira, T. Microbial genome-wide association studies: lessons from human GWAS. *Nat Rev Genet* **18**, 41-50, doi:10.1038/nrg.2016.132 (2017).
- 3 Johnson, R. C. *et al.* Accounting for multiple comparisons in a genome-wide association study (GWAS). *BMC Genomics* **11**, 724, doi:10.1186/1471-2164-11-724 (2010).
- 4 Breiman, L. Statistical Modeling: The Two Cultures. *Statistical Science* **16**, doi:10.1214/ss/1009213726 (2001).
- 5 Read, T. D. & Massey, R. C. Characterizing the genetic basis of bacterial phenotypes using genome-wide association studies: a new direction for bacteriology. *Genome Med* **6**, doi:ARTN 109 10.1186/s13073-014-0109-z (2014).
- 6 Weis, A. M., Clothier, K. A., Huang, B. C., Kong, N. & Weimer, B. C. Draft Genome Sequences of *Campylobacter jejuni* Strains That Cause Abortion in Livestock. *Genome Announc* **4**, doi:10.1128/genomeA.01324-16 (2016).
- 7 Weis, A. M. *et al.* Genomic Comparison of *Campylobacter* spp. and Their Potential for Zoonotic Transmission between Birds, Primates, and Livestock. *Appl Environ Microbiol* **82**, 7165-7175, doi:10.1128/AEM.01746-16 (2016).
- 8 Wu, Z. *et al.* Point mutations in the major outer membrane protein drive hypervirulence of a rapidly expanding clone of *Campylobacter jejuni*. *Proc Natl Acad Sci U S A* **113**, 10690-10695, doi:10.1073/pnas.1605869113 (2016).
- 9 Chen, T. & Guestrin, C. in *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining - KDD '16* 785-794 (2016).
- 10 Behravan, H. *et al.* Machine learning identifies interacting genetic variants contributing to breast cancer risk: A case study in Finnish cases and controls. *Sci Rep* **8**, 13149, doi:10.1038/s41598-018-31573-5 (2018).
- 11 Jason H. Yang, S. N. W., Meagan Hamblin, Douglas McCloskey, Miguel A. Alcantar, Lars Schrübbers, Allison J. Lopatkin, Sangeeta Satish, Amir Nili, Bernhard O. Palsson, Graham C. Walker, James J. Collins. A White-Box Machine Learning Approach for Revealing Antibiotic Mechanisms of Action. *Cell* **177**, 1649-1661, doi:<https://doi.org/10.1016/j.cell.2019.04.016> (2019).
- 12 Heijne, J. N. a. G. Fine-tuning the topology of a polytopic membrane protein: Role of positively and negatively charged amino acids. *Cell* **62**, 1135-1141, doi:[https://doi.org/10.1016/0092-8674\(90\)90390-Z](https://doi.org/10.1016/0092-8674(90)90390-Z) (1990).
- 13 Elazar, A., Weinstein, J. J., Prilusky, J. & Fleishman, S. J. Interplay between hydrophobicity and the positive-inside rule in determining membrane-protein topology. *Proc Natl Acad Sci U S A* **113**, 10340-10345, doi:10.1073/pnas.1605888113 (2016).
- 14 Kim, C. *et al.* Basic amino-acid side chains regulate transmembrane integrin signalling. *Nature* **481**, 209-213, doi:10.1038/nature10697 (2011).
- 15 Seemann, T. Prokka: rapid prokaryotic genome annotation. *Bioinformatics* **30**, 2068-2069, doi:10.1093/bioinformatics/btu153 (2014).

- 332 16 Bandoy, D. Pangenome guided pharmacophore modelling of enterohemorrhagic
333 Escherichia coli sdiA. *F1000Research*
334 doi:<https://doi.org/10.12688/f1000research.17620.1> (2019).
- 335 17 Brown, T. a. l., Luis. sourmash: a library for MinHash sketching of DNA. *Journal of Open*
336 *Source Software* **1**, 27, doi:10.21105/joss.00027 (2016).
- 337 18 Waterhouse, A. *et al.* SWISS-MODEL: homology modelling of protein structures and
338 complexes. *Nucleic Acids Res* **46**, W296-W303, doi:10.1093/nar/gky427 (2018).
- 339 19 Bienert, S. *et al.* The SWISS-MODEL Repository-new features and functionality. *Nucleic*
340 *Acids Res* **45**, D313-D319, doi:10.1093/nar/gkw1132 (2017).
- 341

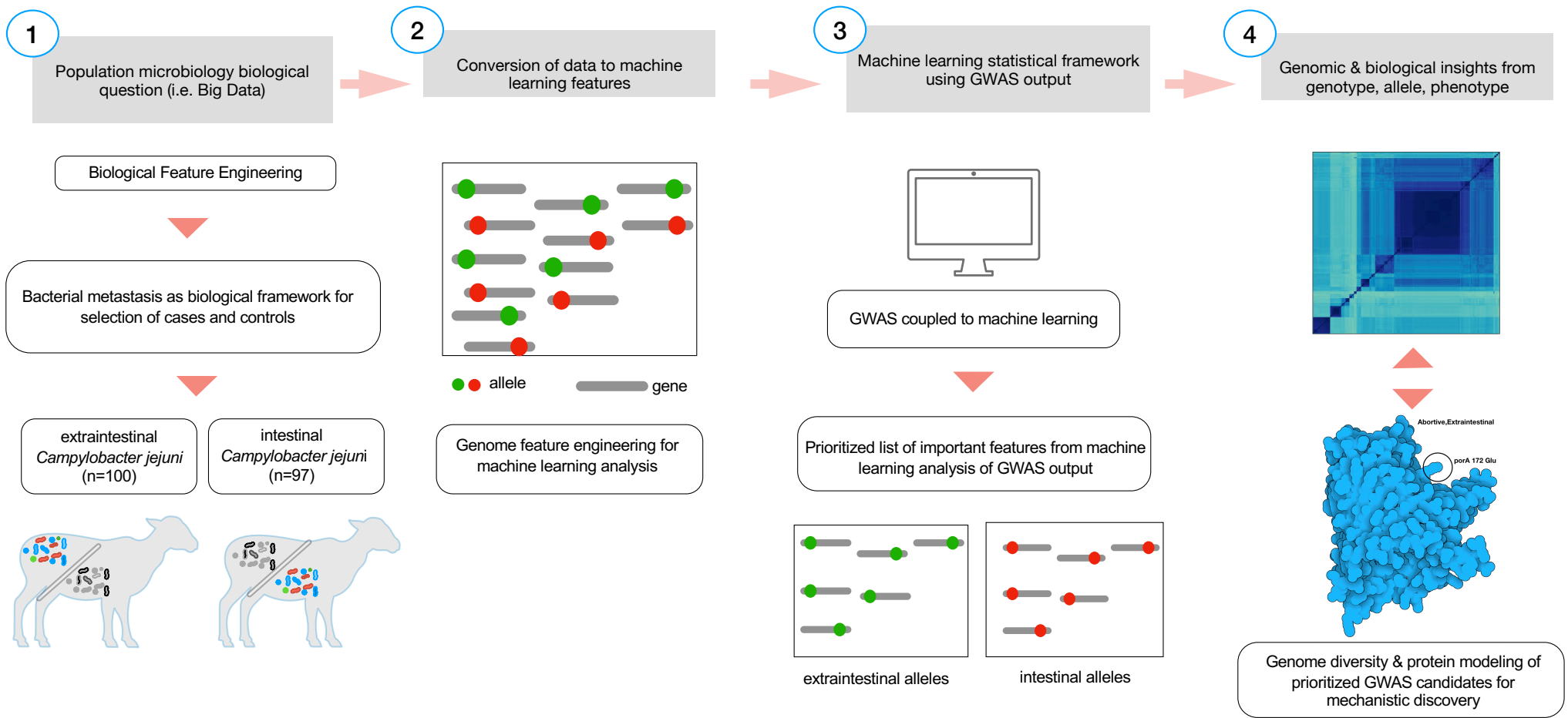


Fig 1.

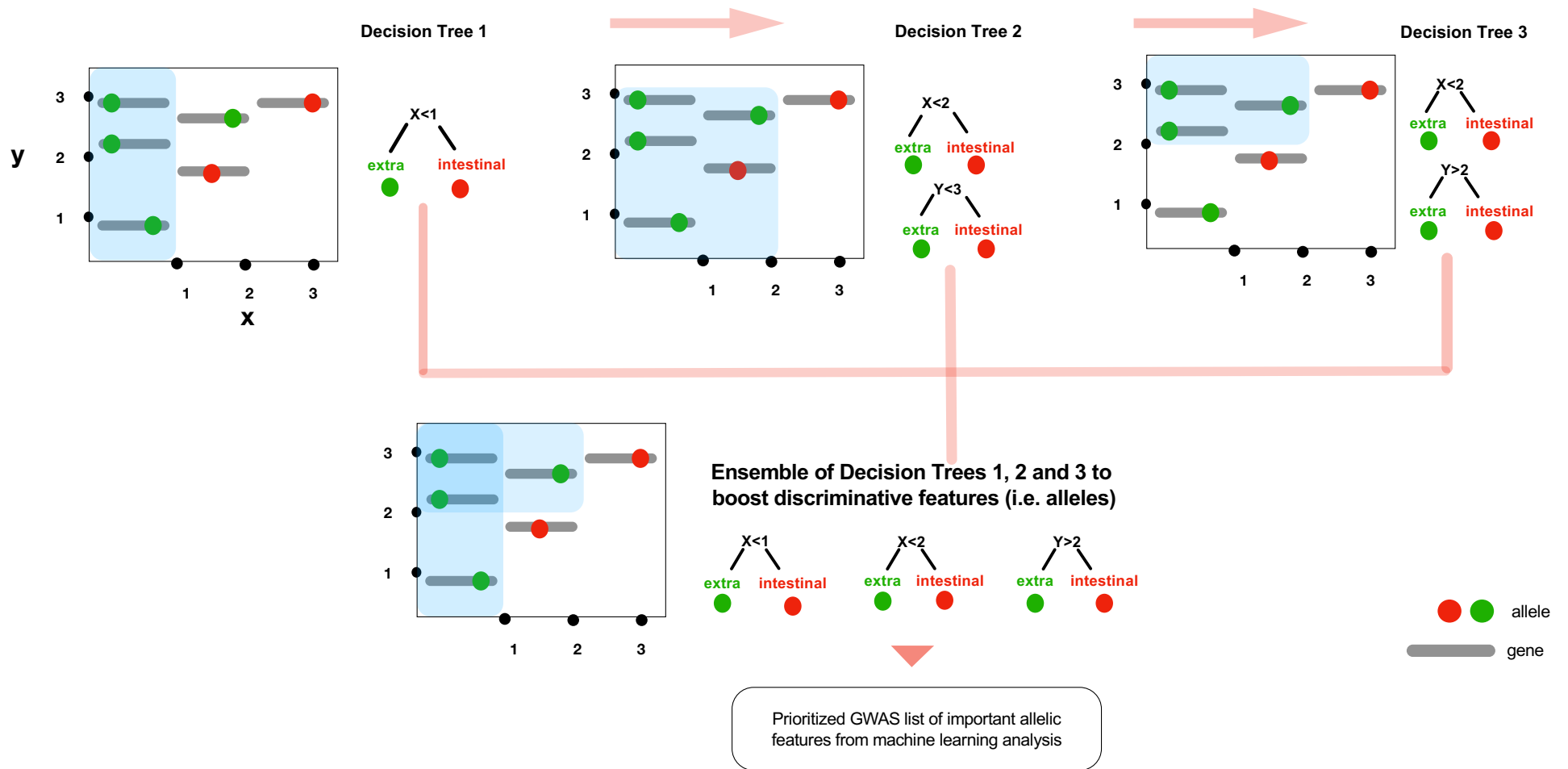


Fig 2.

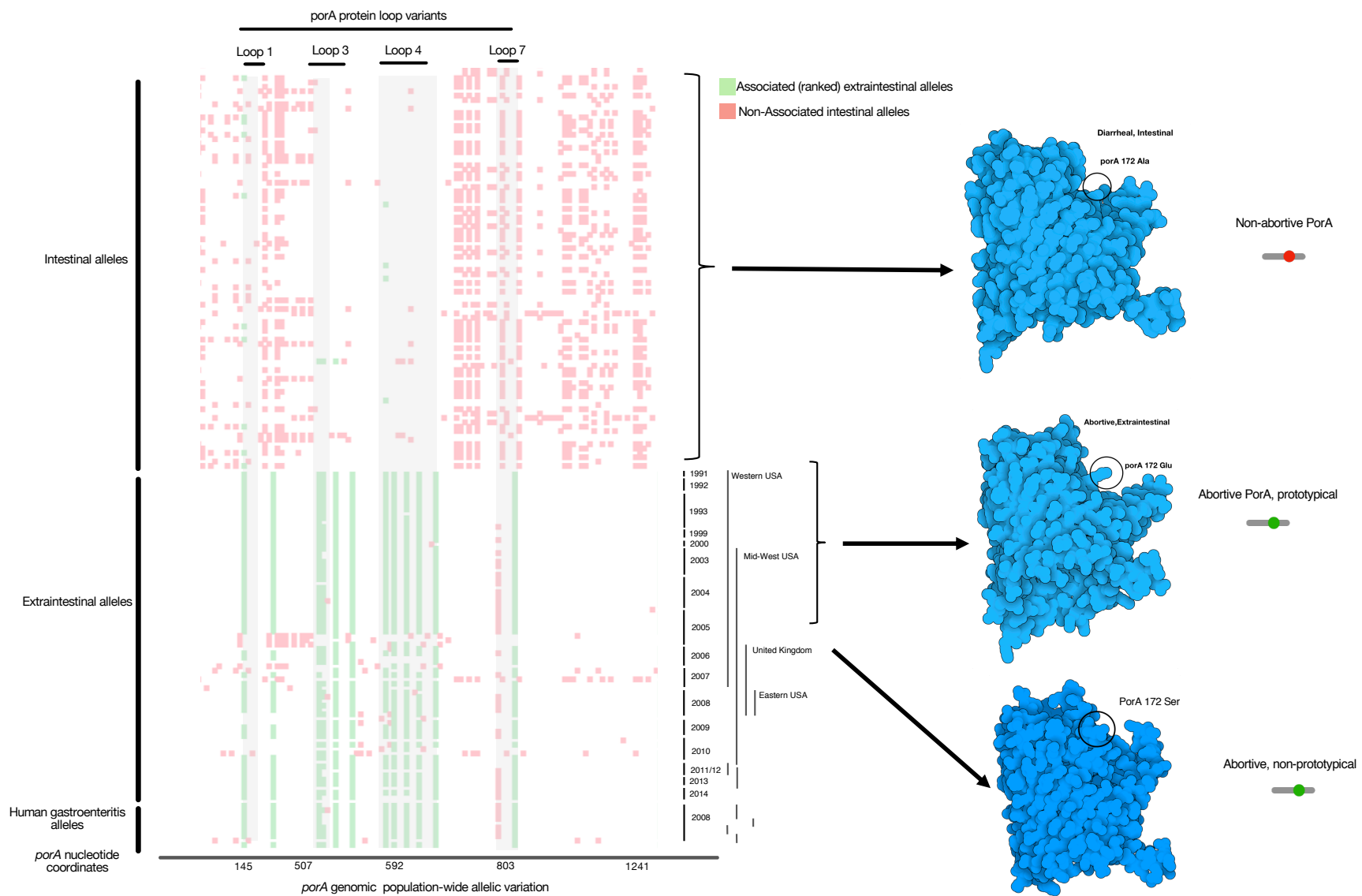


Fig 3.

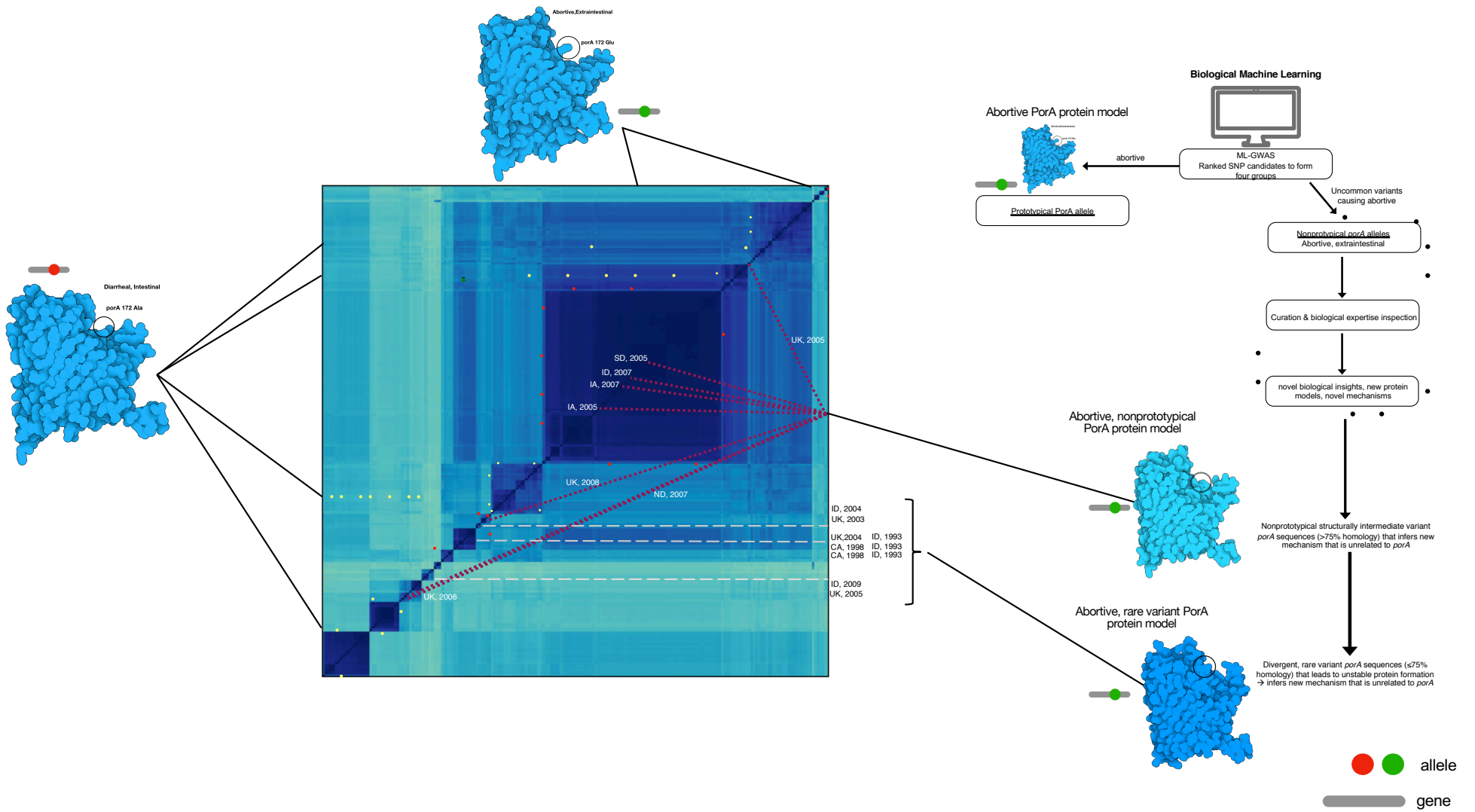
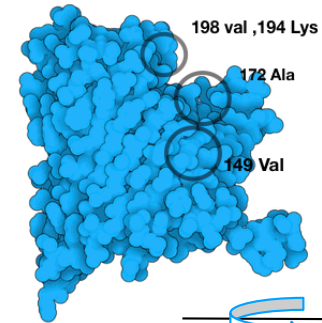


Figure 4

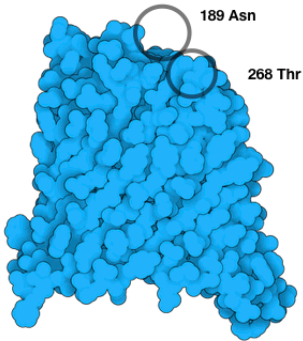
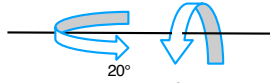
Diarrheal

Intestinal

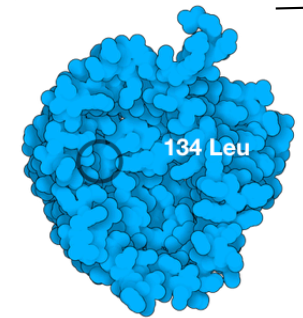
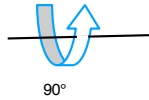


ML-GWAS
Flanked allele with amino acid in the abortive prototypical extraintestinal allele

1. 172 Glu
3. 198 Leu
5. 149 Ile
7. 194 Gln



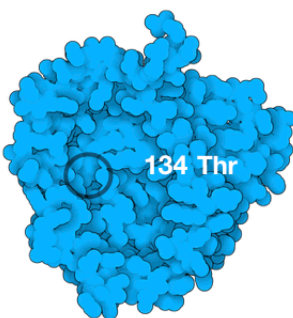
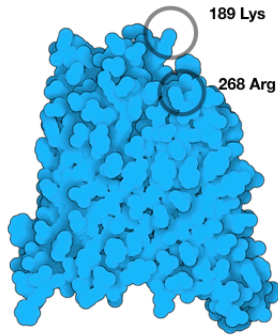
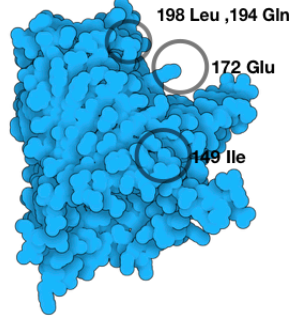
2. 268 Arg
6. 189 Lys



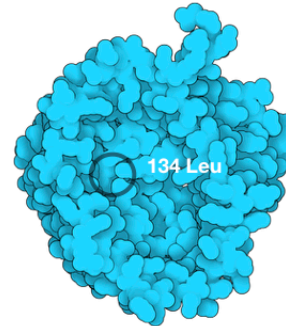
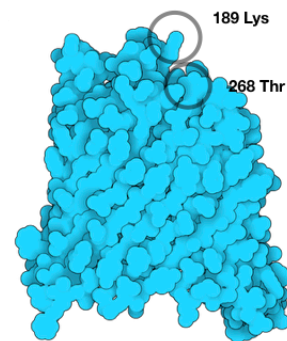
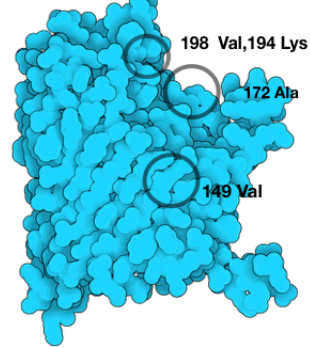
8. 134 Thr

Abortive

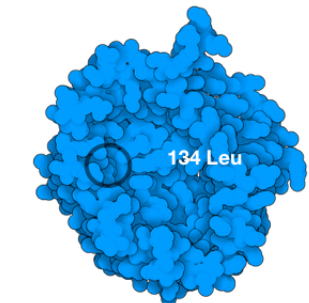
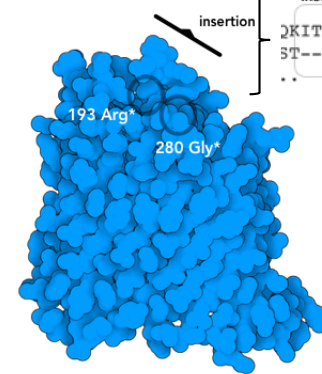
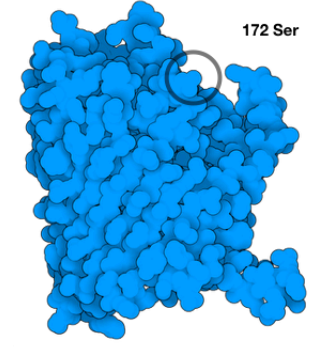
Prototypical extraintestinal



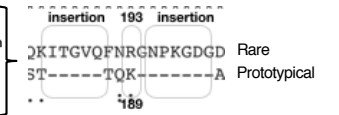
Non-prototypical extraintestinal



Rare variant extraintestinal



counterpart amino acids have been buried deeper, the amino acids are the counterpart of the alleles position as per alignment



The indels resulted in the corresponding amino acids to be buried deeper into the structure. The amino acids were corrected in location number to represent the counterpart of the alleles position as per alignment

Bottom up view of the protein looking up through the barrel that is hidden in this view