

1 **An individualized Bayesian method for estimating genomic variants of hypertension**

2 Md. Asad Rahman<sup>1</sup>, Chunhui Cai<sup>2</sup>, Dennis M. McNamara<sup>3</sup>, Ying Ding<sup>4</sup>, Gregory F. Cooper<sup>2</sup>,  
3 Xinghua Lu<sup>2</sup>, Jinling Liu<sup>1,5\*</sup>

4  
5 <sup>1</sup>Department of Engineering Management and Systems Engineering, Missouri University of  
6 Science and Technology, Rolla, MO, USA, <sup>2</sup>Department of Biomedical Informatics,  
7 University of Pittsburgh, Pittsburgh, PA, USA, <sup>3</sup>Department of Medicine, University of  
8 Pittsburgh, Pittsburgh, PA, USA, <sup>4</sup>Department of Biostatistics, University of Pittsburgh,  
9 Pittsburgh, PA, USA, <sup>5</sup>Department of Biological Sciences, Missouri University of Science and  
10 Technology, Rolla, MO, USA,

11  
12 \*Corresponding author – email: [jinling.liu@mst.edu](mailto:jinling.liu@mst.edu)

## 24 **Abstract**

## 25 **Background**

26           Genomic variants of disease are often discovered nowadays through population-  
27 based genome-wide association studies (GWAS). Identifying genomic variations potentially  
28 underlying a phenotype, such as hypertension, in an individual is important for designing  
29 personalized treatment; however, population-level models, such as GWAS, may not capture  
30 all of the important, individualized factors well. In addition, GWAS typically requires a large  
31 sample size to detect association of low-frequency genomic variants with sufficient power.  
32 Here, we report an individualized Bayesian inference (IBI) algorithm for estimating the  
33 genomic variants that influence complex traits such as hypertension at the level of an  
34 individual (e.g., a patient). By modeling at the level of the individual, IBI seeks to find  
35 genomic variants observed in the individual's genome that provide a strong explanation of  
36 the phenotype observed in this individual.

## 37 **Results**

38           We applied the IBI algorithm to the data from the Framingham Heart Study to explore  
39 genomic influences of hypertension. Among the top-ranking variants identified by IBI and  
40 GWAS, there is a significant number of shared variants (intersection); the unique variants  
41 identified only by IBI tend to have relatively lower minor allele frequency than those identified  
42 by GWAS. In addition, we observed that IBI discovered more individualized and diverse  
43 variants that explain the hypertension patients better than did GWAS. Furthermore, IBI found  
44 several well-known low-frequency variants as well as genes related to blood pressure that  
45 were missed by GWAS in the same cohort. Finally, IBI identified top-ranked variants that  
46 predicted hypertension better than did GWAS, according to the area under the ROC curve.

## 47 **Conclusions**

48           The results provide support for IBI as a promising approach for complementing  
49   GWAS especially in detecting low-frequency genomic variants as well as learning  
50   personalized genomic variants of clinical traits and disease, such as the complex trait of  
51   hypertension, to help advance precision medicine.

52   **Keywords**

53   individualized Bayesian inference, genome-wide association studies, genomic variants,  
54   single nucleotide polymorphism, hypertension, blood pressure, precision medicine

55

56

57

58

59

60

61

62

63

64

65

66

## 67 **Background**

68           Hypertension (HTN) is a key risk factor for many cardiovascular diseases, and it was  
69 primarily responsible for about 7.8 million world-wide deaths in 2015 alone. Previous studies  
70 indicate that in addition to environmental factors, genomic factors play a significant role in  
71 blood pressure (BP) regulation [1]. Hypertension is a polygenic disease [2] burdening a large  
72 population across the globe. Current efforts at identifying significant genomic variants mostly  
73 involve the use of genome-wide association studies (GWAS). Although GWAS has  
74 successfully identified more than 1000 significant single nucleotide polymorphisms (SNPs;  
75 the most common type of genomic variants among people) for BP [3], there are limitations to  
76 this commonly used approach. In general, GWAS requires a large cohort to gain enough  
77 power to identify the significant SNPs, especially the ones with low minor allele frequency  
78 (MAF). That is why before 2015 there were only about 64 significant SNPs identified for  
79 blood pressure, and only recently were more SNPs identified due to the increased sample  
80 sizes (~ 1 million individuals) [4-6]. Still, most of the SNPs identified so far are common  
81 SNPs with small effect sizes, and the total genetic variance in blood pressure explained by  
82 these ~1000 SNPs is small (~5.7%) [3]. It is likely that there are a significant number of non-  
83 common variants missed by GWAS that can help explain much of the remaining genomic  
84 variance [7].

85           GWAS is a population-based approach, and it extracts significant SNPs from a  
86 population level, not considering the specific genome of a given individual. Therefore,  
87 GWAS is not tailored to identify the genomic influences of HTN in an individual, which is the  
88 focus of personalized medicine. It is not uncommon that a HTN patient does not have any of  
89 the significant variants identified at the population level. Thus, identifying the most probable  
90 genomic variants of individual patients is important but remains an unmet need.

91           We have developed an individualized Bayesian inference (IBI) algorithm for  
92 estimating the genomic factors influencing the development of hypertension and other  
93 complex traits in a given individual. As a general machine learning framework, IBI applies a

94 Bayesian method to identify the significant genomic variants in a given individual or patient.  
95 Bayesian methods including Bayesian multiple logistic [8] or linear regression have been  
96 used for identifying the causal SNPs among the significant genomic regions identified by  
97 GWAS (i.e., fine mapping) [9, 10]. However, none of these are individualized. IBI evolved  
98 from a tumor-specific causal inference algorithm (TCI) that members of our team developed  
99 for estimating the somatic mutations driving the development of individual cancerous tumors  
100 [11]. In contrast to TCI, IBI is designed to model and learn the relationships between an  
101 individual genome and a complex trait, such as HTN. Also, IBI was optimized for efficient  
102 computation with whole-genome data, whereas TCI was developed to use whole-exome  
103 data.

104 IBI identifies significant and potentially causal genomic variants for each individual  
105 based on his or her specific genomic background (and available training data on many  
106 similar individuals). By concentrating on the genomic variants observed in a particular  
107 individual, IBI has the potential to discover significant variants of low frequency that exist  
108 only in a small number of individuals and could have been missed by GWAS. The genomic  
109 variants identified being significant by IBI could help inform the design of personalized  
110 treatment for individuals with or at-risk for hypertension.

## 111 **Methods**

### 112 **Overview of Bayesian Networks.**

113 A Bayesian network (BN) [12, 13] is a probabilistic graphical model which has two  
114 components. One is a graphical structure containing nodes and directed edges. Nodes  
115 represent domain variables such as genomic variants or clinical traits. Directed edges  
116 represent conditional dependencies between variables. The other component of a BN is a  
117 set of parameters which are conditional probabilities. Basically, each node has a conditional  
118 probability given its potential causes, which can be described by a conditional probability  
119 function. The joint probabilities of all nodes can be written as a product of each node's

120 conditional probability given its direct causes, based on the local causal Markov condition. A  
121 BN is a flexible framework for modeling the probabilistic relationships among variables in a  
122 complex domain via representing the joint probability of all the variables modeled in a  
123 probabilistic structure. A bipartite BN is a particular class of BN with less complexity, where  
124 there are only two sets of nodes in level 1 and level 2, and potential causal relationships only  
125 occur from nodes in level 1 to nodes in level 2.

126         How do we search for the most probable BN given data? A very popular class of  
127 methods are score-based algorithms that assign a Bayesian score to the BN model and  
128 return the BN with the highest score [12, 13]. This Bayesian score of the BN model is  
129 assigned based on how well this BN is supported by both the data and prior knowledge [14].  
130 In this study, we will use a popular Bayesian score for modeling discrete variables, the  
131 Bayesian Dirichlet equivalent uniform (BDeu) score [14] as TCI did [11].

### 132 **The general framework of individualized Bayesian inference.**

133         As mentioned, IBI is based on TCI [11] and has been further developed and adapted  
134 to fit the circumstances of modeling a variety of complex diseases or traits and whole-  
135 genome genotyping or sequencing data. IBI is designed to estimate the significant genomic  
136 variants, such as SNPs, in a specific individual or patient for downstream clinical and  
137 molecular phenotypes. IBI uses a bipartite BN [12, 13] for modeling the conditional-  
138 dependency or predictive relationships between the genomic variants as a set of  $V$  nodes  
139 and the downstream traits or phenotypes as set of  $T$  nodes; directed edges between  $V$  and  
140  $T$  nodes represent the probabilistic or predictive relationships from variants to traits (Figure  
141 1A). Within this bipartite BN, among all the variants in one individual, IBI assigns a posterior  
142 probability for each variant (represented by an  $V$  node) in influencing or predicting the trait of  
143 interest (represented by node  $T$ ) specific for this individual (Figure 1A, D).

144 For a current individual  $h$ , let  $V_s$  be a variable that represents a specific genomic  
145 variant  $s$  (e.g., a SNP) and let  $T_i$  be a specific trait  $i$  (e.g., HTN) of this individual. Let  $V_s^h$  be a  
146 vector representing all the genomic variants in individual  $h$ . We will examine the predictive  
147 relationship ( $V_s \rightarrow T_i$ ) in this individual for each possible  $V_s$ . Let  $P(V_s \rightarrow T_i)$  be the prior  
148 probability for  $V_s$  predicting or influencing  $T_i$ , which could be estimated using biological  
149 background knowledge or could be set using a uniform prior that assumes all the genomic  
150 variants have the same prior probability of predicting or influencing  $T_i$ . Let  $D$  be the training  
151 data without inclusion of individual  $h$ . Let  $M_s$  represent the log form of the marginal likelihood  
152 of  $P(D | V_s \rightarrow T_i)$  that was derived by assuming one genomic variant  $s$  as the potential  
153 predictor for the entire population  $D$ ;  $M_s$  can be further normalized by the summation of  $M_s$ ,  
154 across all the SNPs to derive the posterior probability (PP). When scoring  $V_s \rightarrow T_i$  for the  
155 entire population  $D$  (i.e., it is not individualized) using Bayesian learning and a uniform prior,  
156 PP is proportional to  $M_s$ ; thus, the ranking of the specific driver  $V_s$  by  $M_s$  or PP as a potential  
157 predictor of a trait at the population level is the same. Thus, we will use  $M_s$  as the score for  
158  $V_s \rightarrow T_i$ . When evaluating the effect of  $V_s$  on  $T_i$  in the entire population using GWAS, the p-  
159 value is derived to indicate the significance of the association between  $V_s$  and  $T_i$  (Figure 1B).

160 IBI partitions the overall population into two subpopulations (Figure 1C). Suppose the  
161 current patient has  $V_s = 1$ , which represents the minor-allele of this SNP. Let  $D^{V_s=1}$  represent  
162 the patient-like-me subpopulation, where all the patients in this subpopulation contain the  
163 value  $V_s = 1$ . IBI evaluates how well  $V_s$  predicts the HTN status within  $D^{V_s=1}$ , which has a  
164 marginal likelihood score of  $P(D^{V_s=1} | V_s \rightarrow T_i)$  that we abbreviate as  $M_s^1$  (Figure 1C, D). Let  
165  $D^{V_s=0}$  represent the remaining cases that do not have  $V_s = 1$ , but rather, have  $V_s = 0$ . To  
166 predict the data in  $D^{V_s=0}$ , IBI finds the SNP  $V_r$  (where “r” denotes the remaining cases) that  
167 maximizes the marginal likelihood of  $V_r \rightarrow T_i$ , namely,  $P(D^{V_s=0} | V_r \rightarrow T_i)$ , which we abbreviate  
168 as  $M_r^0$ . The marginal likelihood for all of the data, given  $V_s$  as an individualized predictor and  
169  $V_r$  as the best predictor of the remaining cases, is  $M_s^1 + M_r^0$ , which we refer to as  $M_{s,r}$

170 (Figure 1C, D). This score of  $M_{s,r}$  can be used to evaluate and rank the capability of  $V_s$  in  
171 explaining the patients-like-me subpopulation that contain this minor allele as well as in  
172 helping reduce the noises for the remaining subpopulation.

173 The marginal likelihood is computed using the BDeu score [14] (Figure 1C, D; refer  
174 to the TCI paper [8]). Individualized posterior probabilities of the form  $V_s \rightarrow T_i$  are further  
175 derived relative to the SNPs that are minor alleles in the genome of the current patient  $h$ .  
176 Thus, the posterior probability takes into consideration the specific genomic background of  
177 the given individual (Figure 1D, Equation 2). In summary, IBI is individualized in the following  
178 ways: (1) The overall marginal likelihood for each arc  $V_s \rightarrow T_i$  (Equation 1) contains an  
179 individualized component that uses the subpopulation of “patients like me” that have the  
180 same variant (i.e.,  $V_s = 1$ ). (2) Each individual has a unique set of genomic variants.  
181 Depending on the specific set of variants, the posterior probability for the same arc of  $V_s \rightarrow T_i$   
182 may be different in different individuals (Equation 2). The individualized nature of IBI makes  
183 it a potential tool for advancing precision medicine where personalized treatments are  
184 desired for individuals of varying genetic backgrounds. IBI is implemented in python with  
185 vectorization and matrix operations for efficient computation involving millions of variants,  
186 and has been tested on whole genome sequencing data on the BioData Catalyst platform  
187 [15].

## 188 **Genome-Wide Association Studies.**

189 GWAS is the standard approach for identifying the significant variants associated  
190 with traits at the population level (e.g., p-value  $< 5 \times 10^{-8}$  for genome wide significance).  
191 Conventional GWAS uses standard logistic regression models or Fisher’s exact test for  
192 discrete traits [16]. We performed GWAS using Fisher’s exact test on the same datasets to  
193 which we applied IBI and compared results.

## 194 **Data and data preprocessing.**

195 We used the whole genome genotyping data of Affymetrix HuGeneFocused50K from  
196 the Framingham Heart Study (FHS) cohort (dbGaP Study Accession: phs000007.v30.p11),



197 which covered about 50K gene-centric and coding SNPs across the genome [17, 18]. We  
198 used the following functions from plink for further filtration and quality control, and have  
199 acquired 38,342 SNPs: --mind 0.03 --geno 0.03 --maf 0.01 --hwe 10e-6 --me 0.05 0.1 --  
200 sexCheck. We filled missing SNP values with the most frequent value for that particular SNP  
201 across the entire population. Dominant coding was then performed in plink, and thus, the  
202 final SNP values are 0 or 1 where 0 represents zero copy of the minor allele (risk allele) and  
203 1 represents one or two copies of the minor allele. The focus of this paper is to predict the  
204 risk (minor) allele SNPs that influence hypertension (high blood pressure) rather than  
205 protecting the subject from hypertension. Therefore, we further removed the SNPs that have  
206 risk ratio smaller than 1 resulting in a total of 19,276 SNPs of interest.

207 Clinical phenotype data included harmonized systolic BP (SBP) and diastolic BP  
208 (DBP) data which were downloaded from PIC-SURE on the NHLBI BioData Catalyst  
209 platform [15]. SBP and DBP are specifically harmonized by the Trans-omics for Precision  
210 Medicine (TOPMed) Data Coordinating Center [19] by taking the average of two SBP or  
211 DBP measurements obtained at a single clinic visit. 10 and 5mm Hg were specifically added  
212 for SBP and DBP for individuals who were taking antihypertensive drugs [20]. If  $SBP \geq 140$ ,  
213 or  $DBP \geq 90$  or an individual who was taking antihypertensive drugs, we considered this  
214 individual as having HTN and assigned 'HTN = 1'; otherwise, we classified this individual as  
215 not having HTN, and we assigned 'HTN = 0'. After merging the SNP data and BP data, we  
216 obtained a total of 6,613 patients with 19,276 SNPs. We performed a stratified random split  
217 to produce an 80% training set (5,290 subjects) and a 20% test set (1,323 subjects), and we  
218 reserved this test set for the prediction task.

## 219 **Results**

220 To evaluate IBI in inferring significant genomic variants for HTN, we compared its  
221 performance to that of GWAS. As a proof of concept, we applied both IBI and GWAS to the  
222 whole genome data of Affymetrix HuGeneFocused50K measurements and harmonized  
223 phenome data of BP measurements from the FHS cohort [18] as described above.

## 224 **A Bayesian method for GWAS analysis.**

225 As was explained in the Methods section, when using a Bayesian method and a  
226 uniform prior to study a single variant's effect on HTN in the population level, the derived  $M_s$   
227 for this variant is proportional to its global posterior probability (GPP), making it possible to  
228 use the marginal likelihood to find the top predicting SNP (Figure 2A). We observed that  
229 when using a uniform prior, as we did in this study, a population-based (i.e., not  
230 individualized) Bayesian approach to identifying top-ranked SNPs based on  $M_s$  yielded  
231 similar results to the population-based GWAS method (Figure 2A). The Spearman  
232 correlation coefficient between the p-value and  $M_s$  across all the 19,276 SNPs is -0.9. We  
233 further examined and compared the top 189 SNPs ranked by  $M_s$  or p-value. The top SNPs  
234 identified by high  $M_s$  values or low p-values are highly overlapping: 164 out of the top 189  
235 SNPs and 18 out of the top 20 SNPs overlapped between these two rankings (Figure 2).  
236 Furthermore, the ranking of these top SNPs by  $M_s$  and p-value are either exactly the same  
237 or very similar where  $M_s$  is negatively correlated with p-value (Figure 2).

## 238 **IBI complements GWAS and better detects significant variants of low MAF.**

239 We applied both IBI and GWAS to the training (discovery) subset of 5,290 FHS  
240 subjects with 19,276 SNPs and HTN status, and derived the IBI marginal values of  $M_{s,r}$   
241 (Figure 3A) and GWAS p-values (Figure 3B) for all the SNPs in the Manhattan plots. In  
242 Figure 3A, the values of  $M_{s,r}$  were normalized with  $(-2436 - M_{s,r}) / (-2436 - (-2463))$   
243 considering -2436 as the maximum and -2463 as the minimum, based on the min-max  
244 normalization technique. For the GWAS analysis, if considering  $0.05 / 19276 = 2.59e-6$  as  
245 the significance level for p-value after the Bonferroni correction, five SNPs reach such  
246 significance (Figure 3B). In Figure 2, the population-level  $M_s$  values were derived by  
247 assuming one SNP as the global predictor or potential cause of HTN for the entire  
248 population (Figure 1B). When using two SNPs to specifically explain HTN status from two  
249 distinct subpopulations as is done by IBI (Figure 1C), the overall marginal values ( $M_{s,r}$ )  
250 significantly increase for many SNPs of  $V_s$ . Among all the SNPs, 189  $V_s$  SNPs have

251  $M_{s,r}$  values bigger than the biggest  $M_s$  value derived in the population level from the best  
252 global predictor, represented as  $V_g$  (Figure 3A). The higher score of IBI compared to the  
253 population-based Bayesian method also has theoretical support. It has been proved that  
254 instance-based (i.e., individualized) causal inference methods, a family of algorithms to  
255 which the IBI belongs, are consistent. More specifically, in the large sample limit, the score  
256 of the data-generating instance-specific model will be assigned the highest score of any  
257 model [21]. These results support that the HTN status in the overall population has been  
258 explained better by IBI with any of the top 189 SNPs,  $V_s$ , explaining the subpopulation of  
259  $D^{V_s=1}$  and with the remaining population predictor,  $V_r$ , explaining the remaining  $D^{V_s=0}$   
260 subpopulation, in comparison to using the best global predictor  $V_g$  itself to explain the entire  
261 population of  $D$ .

262 We performed another evaluation from the perspective of information theory. In this  
263 setting, GWAS analysis is searching for a variant  $V_s$  that has strong information with respect  
264 to a trait  $T_i$  (HTN), and the amount of information can be measured as information gain (IG)  
265 [22]. IG can be calculated by splitting samples according to one variable ( $V_s$ ) and then  
266 measuring the change in the entropy of the other variable ( $T_i$ ) during partitions. Rather than  
267 focusing on just one variant as in a GWAS analysis, the IBI algorithm evaluates how much  
268 information we can gain with respect to the trait  $T_i$  (HTN) if we consider both the specific  
269 variant of interest ( $V_s$ ) and the remaining-population predictor ( $V_r$ ), IG ( $V_s, V_r; T_i$ ) [23]. The  
270 more  $V_s$  and  $V_r$  complement each other (e.g., they are two distinct predictors for different  
271 subgroups) to provide information with respect to  $T_i$ , the higher the IG. In other words, IBI  
272 searches for variants that not only explain HTN well in the “patients-like-me” subpopulation,  
273 but also help enhance the information of  $V_r$  with respect to HTN in the remaining  
274 subpopulation that do not contain this specific variant of interest. Actually, the ranking of the  
275 top 189 SNPs by IBI  $M_{s,r}$  turned out to be highly correlated (Spearman correlation  
276 coefficient,  $r = 0.9$ ) to their ranking by IG values: in general, the higher the  $M_{s,r}$ , the higher  
277 the IG. Based on information gain values, top-5 IBI SNPs were rs11574358, rs1794108,

278 rs11000217, rs2292664 and rs9928967 while top-5 GWAS SNPs were rs11574358,  
279 rs2292664, rs383306, rs5491 and rs1794108. IG values for the top 189 IBI SNPs of  $V_s$   
280 selected by  $M_{s,r}$  are significantly higher than those of the top 189 GWAS SNPs selected by  
281 the p-values; IG values for the top 189 GWAS SNPs are also significantly higher than the  
282 values for 189 randomly-selected SNPs (Figure 3C) as expected.

283 We further examined whether there is any overlap between the top 189 IBI and  
284 GWAS SNPs. We found that 41 of the top 189 SNPs were identified by both IBI and GWAS  
285 (Figure 3D), and 3 of the top 5 IBI and GWAS SNPs are the same (Figure 3A, B). Thus, IBI  
286 and GWAS share many of the same top SNPs, suggesting a mutual agreement between  
287 these two approaches. The unique SNPs for IBI or GWAS support that the two approaches  
288 are also complementary (Fig. 3D). We further examined the MAF distribution for the different  
289 subsets in Figure 3D (Figure 3E). Interestingly, the IBI-only SNPs overall had much lower  
290 MAF than GWAS-only SNPs (Fig. 3E). This result provides support for the hypothesis that  
291 IBI is more capable of identifying lower-frequency significant variants, relative to GWAS, by  
292 concentrating on the genomic variants of a given individual in a specific subpopulation.

293 **IBI discovered more individualized and diverse significant SNPs that better explain**  
294 **the HTN patients, compared to GWAS.**

295 For a given individual  $h$ , IBI derives the posterior probability for each genomic variant  
296  $V_s$ ,  $P(V_s^h \rightarrow T_i^h | D)$ , by normalizing  $M_{s,r}$  with a summation of all the  $M_{s,r}$  across all the existing  
297 minor allele SNPs (i.e., the SNP value is 1) in this individual. This posterior probability takes  
298 into consideration the diverse genomic background or context for different individuals. More  
299 specifically, a particular SNP  $V_s$  with the same  $M_{s,r}$  may have different posterior probabilities  
300 in different individuals due to their distinct genomic background (i.e., different sets of existing  
301 minor allele SNPs). For a given HTN patient, IBI ranked all the minor alleles existing in this  
302 individual based on their individualized posterior probabilities (this ranking will be the same  
303 as the ranking based on  $M_{s,r}$  in a given patient); the SNP with the highest posterior  
304 probability was considered as the most probable influence of HTN for this given patient. For

305 comparison, we designated a top SNP for each HTN patient based on the population-level  
306 p-values derived by GWAS: among the existing minor alleles in a given HTN patient, the  
307 non-protective minor allele with the lowest (most significant) p-value was considered to be  
308 the most probable influence for HTN in this particular patient.

309 Among all the 930 HTN patients in the discovery dataset, we identified 16 unique  
310 SNPs according to GWAS ranking (Figure 4A) and 25 unique SNPs based on IBI ranking  
311 (Figure 4B); each of these unique SNPs was assigned by GWAS or IBI as a top-1 SNP for at  
312 least one HTN patient. The number of HTN patients explained by each of these unique SNP  
313 can be derived by the differences of the accumulated number of explained HTN patients  
314 showed in Figure 4A, B. The more unique SNPs identified by IBI suggested IBI was able to  
315 find a more diverse set of significant SNPs with a more personalized approach. IBI identified  
316 13 SNPs that explain less than 10 HTN patients individually while GWAS found 6 such  
317 SNPs. Interestingly, at the same time, IBI assigns the intronic SNP 'rs13265032' in the  
318 *CSMD1* loci as the top-1 SNP with the highest PP or  $M_{s,r}$  for each of the 425 (46%) of HTN  
319 patients (Figure 4B).

320 For the GWAS analysis, if considering  $0.05/19276 = 2.59e-6$  as the significance level  
321 for p-value after the Bonferroni correction, then 120 out of 930 (12.9%) HTN patients can be  
322 assigned a significant SNP identified by GWAS; even with a relaxed significance level of  
323  $1.09e-5$ , only 146 out of 930 (15.7%) HTN patients are covered or explained by these  
324 significant GWAS SNPs (Figure 4A). This suggests that the significant SNPs of HTN  
325 identified at the population level by GWAS do not necessarily exist in a given HTN patient,  
326 leaving a significant portion of HTN patients unexplained by these significant SNPs.

327 On average, there are 7,767 out of 19,276 minor allele SNPs existing in the HTN  
328 patients. If assuming all these existing risk alleles have the same prior probability in causing  
329 HTN, and only one of them is causing HTN, then the prior probability for each risk allele is  
330  $1.0 / 7767 = 1.3e-4$ . Interestingly, the top one SNP selected by IBI for each HTN patient has  
331 much higher posterior probability ranging from 0.08 to 0.99 (Figure 4C). If considering 0.1 as

332 a significant posterior probability threshold, then 922 out of 930 (99.1%) HTN patients can  
 333 be assigned a significant IBI SNP as the potential cause for their HTN status; with a more  
 334 restrictive threshold of 0.2, 741 out of 930 (79.7%) HTN patients can be explained. These  
 335 results suggested that IBI was able to find a top SNP with significant posterior probability  
 336 ( $\geq 0.1$ ), relative to the random chance ( $1.3e-4$ ), for the majority of HTN patients as a  
 337 potential genomic cause.

338 **Table 1. Novel SNPs in BP-associated genes identified by IBI as the individualized and**  
 339 **most-probable HTN cause**

rs ID	Genes: Variant type	MAF	IBI rank	GWAS rank	p-value	# HTN explained
rs13265032	CSMD1: Intron Variant	0.34	12	5361	0.26	425
rs1564573	CSMD1: Intron Variant	0.42	20	1858	0.08	27
rs2449184	CSMD1: Intron Variant	0.42	33	3760	0.17	0
rs1803274	BCHE: Missense Variant	0.20	23	221	0.01	9
rs948028	GRIK4: Intron Variant	0.16	38	14922	0.78	2
rs12779623	MALRD1: Missense Variant	0.20	48	303	0.01	1

340 As is shown in Table 1, the intronic SNP 'rs13265032' in the *CSMD1* loci that was  
 341 assigned as the top-1 SNP by IBI for 46% (425) of HTN patients, was also ranked high (12)  
 342 by IBI  $M_{s,r}$  among all the SNPs. By contrast, this SNP was never assigned as a top-1 SNP  
 343 for any HTN patient by GWAS and this SNP was ranked very low (5361) by GWAS among  
 344 all the SNPs. Interestingly, other intronic SNPs in the *CSMD1* loci have been reported to be  
 345 associated with hypertension [24] or blood pressure response to hydrochlorothiazide [25,  
 346 26], an antihypertensive drug. Among all the 86 SNPs located in the *CSMD1* loci in our  
 347 dataset, three of them were ranked very high by IBI as is shown in Table 1, while all of them  
 348 were ranked relatively low by GWAS. These three novel SNPs identified by IBI in the  
 349 *CSMD1* loci provide evidence to support the reported role of *CSMD1* in HTN, which  
 350 themselves may warrant further analysis for their potential causal influence on *CSMD1*

351 regulation. In addition to SNPs in the *CSMD1* loci, IBI also identified a novel missense  
352 variant of 'rs1803274' in the *BCHE* loci, a novel intron variant of 'rs948028' in the *GRIK4* loci  
353 and a novel missense variant of 'rs12779623' in the *MALRD1* loci as top-1 likely cause of  
354 HTN in 9, 2 and 1 HTN patient, respectively (Table 1). Interestingly, *BCHE* [27, 28] loci,  
355 *GRIK4* [29] loci and *MALRD1* loci [4] have been reported to be associated with blood  
356 pressure regulation, although GWAS analysis ranked their SNPs relatively low (Table 1).  
357 Overall, these results provide support for IBI being able to identify novel and biologically  
358 meaningful SNPs or genes associated with HTN that were missed by GWAS analysis.

359 **IBI found well-known significant variants or genes that were missed by the parallel**  
360 **GWAS analysis in the same cohort.**

361 We list several missense variants (Table 2) as well as the gene loci (Table 3) that  
362 were previously reported for their influence on blood pressure regulation, in addition to the  
363 ones discussed in Table 1. In Tables 2 and 3, IBI ranks were determined by  $M_{S,r}$  while  
364 GWAS ranks were determined by the p-value.

365 **Table 2. SNPs well-known for blood pressure regulation identified by IBI but missed**  
366 **by GWAS.**

rs ID	Genes: Variant type	MAF	IBI rank	GWAS rank	p-value
rs11575542	DDC: Missense Variant	0.01	79	599	0.02
rs37369	AGXT2: Missense Variant	0.08	93	748	0.03
rs723580	CLIC5: Missense Variant	0.04	189	11851	0.61

367 In Table 2, the missense variant of rs37369 [30] has been shown to be one of the 4  
368 functional SNPs of AGXT2, which has been reported to have strong associations with  
369 several cardiorenal traits, such as coronary heart disease [31]. Its significant association with  
370 hypertension was very recently reported via multiple regression analysis involving only  
371 several targeted SNPs [32]. The missense variant rs11575542 was very recently identified  
372 as a functional variant of the DOPA Decarboxylase (*DDC*) gene during the systematic



373 polymorphism screening across the 15-Exon *DCC* locus [33]. The SNP was shown to alter  
 374 the enzyme activity of DCC and result in changes in renal DA excretion that is linked to  
 375 hypertension [33]. The missense variant of rs723580 was reported to be a top trans-eSNP  
 376 [34] for the expression level of EPO associated with the red blood cell traits that were  
 377 strongly linked to hypertension [35]. Intriguingly, with low MAF in our relatively small  
 378 discovery cohort, these three SNPs were ranked much higher by IBI than by the parallel  
 379 GWAS analysis (Table 2). This result provides support that IBI can recognize biologically  
 380 meaningful genomic variants of low MAF, relative to GWAS, particularly when the sample  
 381 size is small compared to the number of SNPs to be tested.

382 **Table 3. Genes well-known for blood pressure regulation identified by IBI but ranked**  
 383 **relatively low by GWAS.**

rs ID	Genes: Variant type	MAF	IBI rank	GWAS rank	Top GWAS rank	p-value
rs3211938	CD36 [36, 37]: Stop Gained	0.01	32	141	141	3.8E-03
rs6730396	ALLC [38]: Missense Variant	0.01	45	192	192	5.5E-03
rs9896904	ANKFN1 [39]: Intron Variant	0.07	57	14392	2130	7.5E-01
rs11899922	THSD7B [4, 40]: Intron Variant	0.07	70	9415	929	4.8E-01
rs10968668	LINGO2 [41, 42]: Intron Variant	0.08	80	1237	1237	4.9E-02
rs13261739	PDGFRL [43, 44]: Intron Variant	0.13	94	7156	7156	3.6E-01
rs6140644	PLCB1 [45]: Intron Variant	0.17	116	9947	699	5.1E-01
rs7647302	KCNAB1 [46]: Intron Variant	0.01	158	9528	1964	4.9E-01

384 Table 3 shows a list of genes that have been reported to be associated with BP  
 385 regulation or HTN where at least one related paper is listed for each gene. IBI identified  
 386 these genes as candidate genes influencing HTN since these gene loci contain at least one  
 387 novel SNP that is highly ranked by IBI for its association with HTN. Interestingly, all of these



388 genes were ranked relatively low by GWAS even considering the highest SNP rank by  
389 GWAS ('Top GWAS rank' in Table 3) within each gene locus; all of these novel SNPs were  
390 also ranked relatively low by GWAS ('GWAS rank' in Table 3) with non-significant p-values  
391 (Table 3). Moreover, among these eight SNPs highly-ranked by IBI and lowly-ranked by  
392 GWAS, six have MAF lower than 0.1 and three have MAF as low as 0.01 (Table 3). This  
393 together with Table 2 supports that IBI is more capable in identifying significant variants of  
394 low MAF, compared to GWAS.

395 **IBI top SNPs identify genetic risk scores that are more predictive for HTN than do the**  
396 **GWAS top SNPs.**

397 We further compared the capabilities of significant SNPs, identified by IBI and  
398 GWAS, in predicting HTN. After running IBI and GWAS on the training (discovery) dataset,  
399 we were able to rank all the SNPs based on  $M_{s,r}$  derived from IBI or p-values obtained from  
400 GWAS. For each subject in the test set, based on IBI ranking or GWAS ranking, we  
401 identified the top 1 and 3 SNPs that exist in this subject (with a value of '1' denoting a minor  
402 allele). We then used these top SNPs to calculate the genetic risk scores (GRS) for each  
403 subject by the sum of his or her risk alleles, weighted by odds ratio for GWAS top SNPs or  
404 by  $M_{s,r}$  for IBI top SNPs. We further use min-max normalization to normalize both the IBI  
405 and GWAS GRS to avoid potential bias from the different scales of the original values. We  
406 then directly calculated the area under ROC curve (AUROC or AUC) using the normalized  
407 GRS for each patient (Figure 5). We also trained a logistic regression model for HTN  
408 prediction using this feature of normalized GRS, which gave very similar results (data not  
409 shown).

410 As expected, using randomly selected SNPs showed poor prediction performance,  
411 with an AUROC of 0.50 (Figure 5A, D); the GWAS-selected top one or three SNPs both  
412 have an AUROC of 0.55 (Figure 5B, E) suggesting some level of prediction; the IBI-selected  
413 top SNP had an AUROC of 0.59 (Figure 5C) and the top 3 SNPs had an AUROC of 0.60

414 (Figure 5F). The higher AUROC achieved by IBI provides support that the top SNPs it  
415 selects predict hypertension better than the top SNPs selected by GWAS.

## 416 **Discussion**

417 In this study, we developed and applied a novel and individualized method (IBI) to  
418 estimate the personalized genomic variants for the complex trait of hypertension. We  
419 compared its performance with the population-based GWAS method using a real dataset  
420 from the FHS cohort. The significant overlap of the top-ranked SNPs by both IBI and GWAS  
421 suggest a degree of agreement of these two approaches. On the other hand, the unique  
422 SNPs they found support a complementary role of IBI to current GWAS analyses.  
423 Interestingly, by focusing on each individual and its patient-like-me subgroup, IBI was  
424 capable of identifying significant SNPs of low MAF in the same cohort, relative to GWAS. IBI  
425 was also able to identify more diverse and individualized top SNPs to explain the HTN  
426 patients. Moreover, the top SNPs identified by IBI from the discovery cohort were able to  
427 predict HTN better than the top ones derived from GWAS when applied to an unseen test  
428 cohort. We also identified evidence from the literature to support the biological significance  
429 of top SNPs found by IBI, especially the ones highly-ranked by IBI and lowly-ranked by  
430 GWAS. In summary, our study provides support that IBI can serve as a complementary  
431 approach in discovering novel and personalized genomic variants that may be missed by  
432 GWAS.

433 Contemporary GWAS studies often involve a large sample size (~1 million) to gain  
434 sufficient power, especially for variants of low MAF. Considering the large genomic  
435 heterogeneity among different individuals, as well as the nature of complex diseases often  
436 being affected by many variants of small effect size, an alternative approach is to focus on  
437 the subpopulation containing the specific variant of low MAF under evaluation, as IBI does.  
438 In this way, IBI may be able to better evaluate the effect of low-MAF variants in a patient-  
439 like-me subpopulation, without the potential noise from a large remaining population not

440 containing such variants; moreover, this large remaining population could be explained  
441 better with a remaining population driver. The fact that the top-ranked SNPs by IBI in general  
442 have a higher overall marginal likelihood,  $M_{s,r}$ , and higher information gain with respect to  
443 the HTN status, provide support that IBI may have found specific drivers that better explain  
444 the subpopulations. Our results also support that IBI is not compromised in identifying  
445 significant high-MAF SNPs. IBI's population partition strategy aligns well with the concept of  
446 personalized medicine in which different individuals or subpopulations may have different  
447 underlying genomic influences on producing complex clinical phenotypes such as HTN.

448         As a general Bayesian framework, IBI can be applied to any discrete trait. It can also  
449 be applied to continuous traits by changing the marginal likelihood function from using the  
450 BDeu score for discrete variables to using the Bayesian information criterion (BIC) score for  
451 continuous variables [14]. For the current approach presented in this study, one limitation is  
452 that it only considers the genomic factors of HTN, while not modeling the effects of other  
453 factors such as age, sex, population structure and the family relatedness that may exist in  
454 this FHS cohort. To model the effects from non-genomic factors, we plan to incorporate  
455 linear mixed models [47-51] into our current framework. Also, due to confounding factors  
456 such as population structure, as well as linkage disequilibrium (LD), the predictive variants  
457 described in this paper are not guaranteed to be causal. Further fine mapping approaches,  
458 functional analysis, or Mendelian randomization can be used to further pinpoint the potential  
459 causality. Another interesting future direction is to search for more than one genomic variant  
460 that might work together to affect and predict the phenotypes of individuals and  
461 subpopulations.

## 462 **Conclusions**

463         In summary, we described a novel Bayesian method for identifying personalized  
464 genomic variants that predict complex traits, such as HTN. IBI can serve as a  
465 complementary approach to GWAS, especially in detecting significant genomic variants of

466 low frequency. The novel SNPs we identified for HTN warrant further analysis for their  
467 possible causal role in blood pressure regulation.

#### 468 **List of abbreviations**

469 GWAS: Genome Wide Association Study

470 IBI: Individualized Bayesian inference

471 TCI: tumor-specific causal inference

472 HTN: Hypertension

473 BP: Blood pressure

474 SBP: Systolic blood pressure

475 DBP: Diastolic blood pressure

476 FHS: Framingham Heart Study

477 MAF: Minor allele frequency

478 CBN: Causal Bayesian network

479 BDeu: Bayesian Dirichlet equivalent uniform

480 TOPMed: The Trans-Omics for Precision Medicine

481 IG: Information gain

482 GRS: Genetic risk score

483 AUROC or AUC: area under ROC curve

484 BIC: Bayesian information criterion

#### 485 **Declarations**

#### 486 **Ethics approval and consent to participate**

487 Not applicable.

#### 488 **Consent for publication**

489 Not applicable.

#### 490 **Availability of data and materials**

491 The genomic and blood pressure data from the FHS cohort (study  
492 accession: phs000007.v30.p11) are publicly available via controlled-access through the  
493 database of Genotypes and Phenotypes Study  
494 ([https://www.ncbi.nlm.nih.gov/projects/gap/cgi-bin/study.cgi?study\\_id=phs000007.v30.p11](https://www.ncbi.nlm.nih.gov/projects/gap/cgi-bin/study.cgi?study_id=phs000007.v30.p11)).

495 The source code for IBI is publicly available on GitHub at <https://github.com/asadcfc/IBI>.

#### 496 **Competing interests**

497 The authors declared no competing interests.

#### 498 **Authors' contributions**

499 JL designed and implemented the IBI algorithm in python, conceived and designed the  
500 experiments, and wrote the paper; MA.R and JL carried out the data collection, modeling,  
501 analyses; XL and GC provided insightful advices about the design of the IBI algorithm and  
502 some of the experiments, as well as edited the manuscript; YD provided advice for whole-  
503 genome data processing and GWAS analysis; DM provided general advice on hypertension  
504 development and its potential genomic causes; CC helped with the information gain  
505 experiment; All authors read and approved the final manuscript.

#### 506 **Acknowledgements**

507 The Framingham Heart Study, within the database of Genotypes and Phenotypes Study  
508 (accession #: phs000007.v30.p11), is conducted and supported by the National Heart, Lung,  
509 and Blood Institute (NHLBI) in collaboration with Boston University (Contract No. N01-HC-  
510 25195 and HHSN268201500001I). This manuscript was not prepared in collaboration with  
511 investigators of the Framingham Heart Study and does not necessarily reflect the opinions  
512 or views of the Framingham Heart Study, Boston University, or NHLBI. Funding for SHARe  
513 Affymetrix genotyping was provided by NHLBI Contract N02-HL- 64278. SHARe Illumina  
514 genotyping was provided under an agreement between Illumina and Boston University.

515 Funding for Affymetrix genotyping of the FHS Omni cohorts was provided by Intramural  
516 NHLBI funds from Andrew D. Johnson and Christopher J. O'Donnell.

517 Support for this work was provided by the National Institutes of Health, National Heart, Lung,  
518 and Blood Institute, through the BioData Catalyst program (award 1OT3HL142479-01,  
519 1OT3HL142478-01, 1OT3HL142481-01, 1OT3HL142480-01, 1OT3HL147154-01). Any  
520 opinions expressed in this document are those of the author(s) and do not necessarily reflect  
521 the views of NHLBI, individual BioData Catalyst team members, or affiliated organizations  
522 and institutions. The authors wish to acknowledge the contributions of the consortium  
523 working on the development of the NHLBI BioData Catalyst ecosystem.

524

## 525 **Funding**

526 This research was funded in part by the BioData Catalyst Fellowship award (0065304) from  
527 the National Institutes of Health, National Heart, Lung, and Blood Institute through the  
528 University of North Carolina at Chapel Hill, grant R01LM012011 from the National Institutes  
529 of Health, R01 MD009118 from the National Institute on Minority Health and Health  
530 Disparities (NIMHD).

## 531 **References**

- 532 1. Padmanabhan S, Joe B: Towards Precision Medicine for Hypertension: A Review of  
533 Genomic, Epigenomic, and Microbiomic Effects on Blood Pressure in Experimental Rat  
534 Models and Humans. *Physiol Rev* 2017, 97(4):1469-1528.
- 535 2. Padmanabhan S, Dominiczak AF: Genomics of hypertension: the road to precision  
536 medicine. *Nat Rev Cardiol* 2021, 18(4):235-250.
- 537 3. Cabrera CP, Ng FL, Nicholls HL, Gupta A, Barnes MR, Munroe PB, Caulfield MJ: Over  
538 1000 genetic loci influencing blood pressure with multiple systems and tissues  
539 implicated. *Hum Mol Genet* 2019, 28(R2):R151-R161.

- 540 4. Evangelou E, Warren HR, Mosen-Ansorena D, Mifsud B, Pazoki R, Gao H, Ntritsos G,  
541 Dimou N, Cabrera CP, Karaman I *et al*: Genetic analysis of over 1 million people  
542 identifies 535 new loci associated with blood pressure traits. *Nat Genet* 2018,  
543 50(10):1412-1425.
- 544 5. Surendran P, Feofanova EV, Lahrouchi N, Ntalla I, Karthikeyan S, Cook J, Chen L,  
545 Mifsud B, Yao C, Kraja AT *et al*: Discovery of rare variants associated with blood  
546 pressure regulation through meta-analysis of 1.3 million individuals. *Nat Genet* 2020,  
547 52(12):1314-1332.
- 548 6. Giri A, Hellwege JN, Keaton JM, Park J, Qiu C, Warren HR, Torstenson ES, Kovesdy  
549 CP, Sun YV, Wilson OD *et al*: Trans-ethnic association study of blood pressure  
550 determinants in over 750,000 individuals. *Nat Genet* 2019, 51(1):51-62.
- 551 7. Russo A, Di Gaetano C, Cugliari G, Matullo G: Advances in the Genetics of  
552 Hypertension: The Effect of Rare Variants. *Int J Mol Sci* 2018, 19(3).
- 553 8. Banerjee S, Zeng L, Schunkert H, Soding J: Bayesian multiple logistic regression for  
554 case-control GWAS. *PLoS Genet* 2018, 14(12):e1007856.
- 555 9. Schaid DJ, Chen W, Larson NB: From genome-wide associations to candidate causal  
556 variants by statistical fine-mapping. *Nat Rev Genet* 2018, 19(8):491-504.
- 557 10. Stephens M, Balding DJ: Bayesian statistical methods for genetic association studies.  
558 *Nat Rev Genet* 2009, 10(10):681-690.
- 559 11. Cai C, Cooper GF, Lu KN, Ma X, Xu S, Zhao Z, Chen X, Xue Y, Lee AV, Clark N *et al*:  
560 Systematic discovery of the functional impact of somatic genome alterations in individual  
561 tumors through tumor-specific causal inference. *PLoS Comput Biol* 2019,  
562 15(7):e1007088.
- 563 12. Heckerman D, Meek C, Cooper G: A Bayesian approach to causal discovery. In:  
564 *Computation, causation, and discovery (MIT Press)* Edited by Glymour C, Cooper G,  
565 vol. 19; 1999: 141-166.
- 566 13. Spirtes P GC, Scheines R.: *Causation, Prediction, and Search.*: Cambridge, MA: MIT  
567 Press;; 2000.

- 568 14. Heckerman D, Geiger D, Chickering DM: Learning Bayesian networks: The combination  
569 of knowledge and statistical data. *Machine learning* 1995, 20(3):197-243.
- 570 15. Consortium BC: The NHLBI BioData Catalyst. *Zenodo*(31 August 2020, date last  
571 accessed) 2020.
- 572 16. Gogarten SM, Sofer T, Chen H, Yu C, Brody JA, Thornton TA, Rice KM, Conomos MP:  
573 Genetic association testing using the GENESIS R/Bioconductor package. *Bioinformatics*  
574 2019, 35(24):5346-5348.
- 575 17. Mahmood SS, Levy D, Vasani RS, Wang TJ: The Framingham Heart Study and the  
576 epidemiology of cardiovascular disease: a historical perspective. *The Lancet* 2014,  
577 383(9921):999-1008.
- 578 18. Tsao CW, Vasani RS: Cohort Profile: The Framingham Heart Study (FHS): overview of  
579 milestones in cardiovascular epidemiology. *Int J Epidemiol* 2015, 44(6):1800-1813.
- 580 19. Stilp AM, Emery LS, Broome JG, Buth EJ, Khan AT, Laurie CA, Wang FF, Wong Q,  
581 Chen D, D'Augustine CM *et al*: A System for Phenotype Harmonization in the National  
582 Heart, Lung, and Blood Institute Trans-Omics for Precision Medicine (TOPMed)  
583 Program. *Am J Epidemiol* 2021, 190(10):1977-1992.
- 584 20. Levy D, Ehret GB, Rice K, Verwoert GC, Launer LJ, Dehghan A, Glazer NL, Morrison  
585 AC, Johnson AD, Aspelund T *et al*: Genome-wide association study of blood pressure  
586 and hypertension. *Nat Genet* 2009, 41(6):677-687.
- 587 21. Jabbari F: Instance-Specific Causal Bayesian Network Structure Learning. University of  
588 Pittsburgh; 2021.
- 589 22. KENT JT: Information gain and a general measure of correlation. *Biometrika* 1983,  
590 70(1):163-173.
- 591 23. Quinlan JR: Induction of decision trees. *Machine learning* 1986, 1(1):81-106.
- 592 24. Hong KW, Go MJ, Jin HS, Lim JE, Lee JY, Han BG, Hwang SY, Lee SH, Park HK, Cho  
593 YS *et al*: Genetic variations in ATP2B1, CSK, ARSG and CSMD1 loci are related to  
594 blood pressure and/or hypertension in two Korean cohorts. *Journal of Human*  
595 *Hypertension* 2010, 24(6):367-372.



- 596 25. Chittani M, Zaninello R, Lanzani C, Frau F, Ortu MF, Salvi E, Fresu G, Citterio L, Braga  
597 D, Piras DA *et al*: TET2 and CSMD1 genes affect SBP response to hydrochlorothiazide  
598 in never-treated essential hypertensives. *J Hypertens* 2015, 33(6):1301-1309.
- 599 26. Salvi E, Wang Z, Rizzi F, Gong Y, McDonough CW, Padmanabhan S, Hiltunen TP,  
600 Lanzani C, Zaninello R, Chittani M *et al*: Genome-Wide and Gene-Based Meta-Analyses  
601 Identify Novel Loci Influencing Blood Pressure Response to Hydrochlorothiazide.  
602 *Hypertension* 2017, 69(1):51-59.
- 603 27. Benyamin B, Middelberg RP, Lind PA, Valle AM, Gordon S, Nyholt DR, Medland SE,  
604 Henders AK, Heath AC, Madden PAF *et al*: GWAS of butyrylcholinesterase activity  
605 identifies four novel loci, independent effects within BCHE and secondary associations  
606 with metabolic risk factors. *Human Molecular Genetics* 2011, 20(22):4504-4514.
- 607 28. Cardoso AM, Abdalla FH, Bagatini MD, Martins CC, Fiorin Fda S, Baldissarelli J, Costa  
608 P, Mello FF, Fiorenza AM, Serres JD *et al*: Swimming training prevents alterations in  
609 acetylcholinesterase and butyrylcholinesterase activities in hypertensive rats. *Am J*  
610 *Hypertens* 2014, 27(4):522-529.
- 611 29. Rutherford S, Cai G, Lopez-Alvarenga JC, Kent JW, Voruganti VS, Proffitt JM, Curran  
612 JE, Johnson MP, Dyer TD, Jowett JB: A chromosome 11q quantitative-trait locus  
613 influences change of blood-pressure measurements over time in Mexican Americans of  
614 the San Antonio Family Heart Study. *The American Journal of Human Genetics* 2007,  
615 81(4):744-755.
- 616 30. Suhre K, Wallaschofski H, Raffler J, Friedrich N, Haring R, Michael K, Wasner C, Krebs  
617 A, Kronenberg F, Chang D: A genome-wide association study of metabolic traits in  
618 human urine. *Nature genetics* 2011, 43(6):565-569.
- 619 31. Hu XL, Zeng WJ, Li MP, Yang YL, Kuang DB, Li H, Zhang YJ, Jiang C, Peng LM, Qi H  
620 *et al*: AGXT2 rs37369 polymorphism predicts the renal function in patients with chronic  
621 heart failure. *Gene* 2017, 637:145-151.
- 622 32. Yoshino Y, Kumon H, Mori T, Yoshida T, Tachibana A, Shimizu H, Iga J-i, Ueno S-i:  
623 Effects of AGXT2 variants on blood pressure and blood sugar among 750 older

- 624 Japanese subjects recruited by the complete enumeration survey method. *BMC*  
625 *Genomics* 2021, 22(1):287.
- 626 33. Miramontes-Gonzalez JP, Hightower CM, Zhang K, Kurosaki H, Schork AJ, Biswas N,  
627 Vaingankar S, Mahata M, Lipkowitz MS, Nievergelt CM *et al*: A new common functional  
628 coding variant at the DDC gene change renal enzyme activity and modify renal  
629 dopamine function. *Scientific Reports* 2019, 9(1):5055.
- 630 34. Zhang X, Johnson AD, Hendricks AE, Hwang S-J, Tanriverdi K, Ganesh SK, Smith NL,  
631 Peyser PA, Freedman JE, O'Donnell CJ: Genetic associations with expression for  
632 genes implicated in GWAS studies for atherosclerotic cardiovascular disease and blood  
633 phenotypes. *Human Molecular Genetics* 2013, 23(3):782-795.
- 634 35. Tsuda K: Red blood cell abnormalities and hypertension. *Hypertension Research* 2020,  
635 43(1):72-73.
- 636 36. Momeni-Moghaddam MA, Asadikaram G, Akbari H, Abolhassani M, Masoumi M,  
637 Nadimy Z, Khaksari M: CD36 gene polymorphism rs1761667 (G> A) is associated with  
638 hypertension and coronary artery disease in an Iranian population. *BMC cardiovascular*  
639 *disorders* 2019, 19(1):1-9.
- 640 37. Pravenec M, Churchill PC, Churchill MC, Viklicky O, Kazdova L, Aitman TJ, Petretto E,  
641 Hubner N, Wallace CA, Zimdahl H: Identification of renal Cd36 as a determinant of  
642 blood pressure and risk for hypertension. *Nature genetics* 2008, 40(8):952-954.
- 643 38. Sung YJ, Basson J, Cheng N, Nguyen K-DH, Nandakumar P, Hunt SC, Arnett DK,  
644 Dávila-Román VG, Rao DC, Chakravarti A: The role of rare variants in systolic blood  
645 pressure: analysis of ExomeChip data in HyperGEN African Americans. *Human heredity*  
646 2015, 79(1):20-27.
- 647 39. Chittani M, Zaninello R, Lanzani C, Frau F, Ortu MF, Salvi E, Fresu G, Citterio L, Braga  
648 D, Piras DA: TET2 and CSMD1 genes affect SBP response to hydrochlorothiazide in  
649 never-treated essential hypertensives. *Journal of hypertension* 2015, 33(6):1301.

- 650 40. de Las Fuentes L, Sung YJ, Schwander KL, Kalathiveetil S, Hunt SC, Arnett DK, Rao D:  
651 The role of SNP-loop diuretic interactions in hypertension across ethnic groups in  
652 HyperGEN. *Frontiers in genetics* 2013, 4:304.
- 653 41. Parmar PG, Taal HR, Timpson NJ, Thiering E, Lehtimäki T, Marinelli M, Lind PA, Howe  
654 LD, Verwoert G, Aalto V: International genome-wide association study consortium  
655 identifies novel loci associated with blood pressure in children and adolescents.  
656 *Circulation: Cardiovascular Genetics* 2016, 9(3):266-278.
- 657 42. Feitosa MF, Kraja AT, Chasman DI, Sung YJ, Winkler TW, Ntalla I, Guo X, Franceschini  
658 N, Cheng C-Y, Sim X: Novel genetic associations for blood pressure identified via gene-  
659 alcohol interaction in up to 570K individuals across multiple ancestries. *PloS one* 2018,  
660 13(6):e0198166.
- 661 43. Yang W, Huang J, Yao C, Fan Z, Ge D, Gan W, Huang G, Hui R, Shen Y, Qiang B:  
662 Evidence for linkage and association of the markers near the LPL gene with  
663 hypertension in Chinese families. *Journal of medical genetics* 2003, 40(5):e57-e57.
- 664 44. Chun HJ, Bonnet S, Chan SY: Translational advances in the field of pulmonary  
665 hypertension. Translating MicroRNA biology in pulmonary hypertension. It will take more  
666 than “miR” words. *American journal of respiratory and critical care medicine* 2017,  
667 195(2):167-178.
- 668 45. Warren HR, Evangelou E, Cabrera CP, Gao H, Ren M, Mifsud B, Ntalla I, Surendran P,  
669 Liu C, Cook JP: Genome-wide association analysis identifies novel blood pressure loci  
670 and offers biological insights into cardiovascular risk. *Nature genetics* 2017, 49(3):403-  
671 415.
- 672 46. McCarthy NS, Vangjeli C, Cavalleri GL, Delanty N, Shianna KV, Surendran P, O'Brien  
673 E, Munroe PB, Masca N, Tomaszewski M: Two further blood pressure loci identified in  
674 ion channel genes with a gene-centric approach. *Circulation: Cardiovascular Genetics*  
675 2014, 7(6):873-879.
- 676 47. Chen H, Wang C, Conomos MP, Stilp AM, Li Z, Sofer T, Szpiro AA, Chen W, Brehm JM,  
677 Celedon JC *et al*: Control for population structure and relatedness for binary traits in

678 genetic association studies via logistic mixed models. *Am J Hum Genet* 2016,  
679 98(4):653-666.

680 48. Lippert C, Listgarten J, Liu Y, Kadie CM, Davidson RI, Heckerman D: FaST linear mixed  
681 models for genome-wide association studies. *Nat Methods* 2011, 8(10):833-835.

682 49. Pirinen M, Donnelly P, Spencer CCA: Efficient computation with a linear mixed model on  
683 large-scale data sets with applications to genetic studies. *The Annals of Applied*  
684 *Statistics* 2013, 7(1).

685 50. Wen YJ, Zhang H, Ni YL, Huang B, Zhang J, Feng JY, Wang SB, Dunwell JM, Zhang  
686 YM, Wu R: Methodological implementation of mixed linear models in multi-locus  
687 genome-wide association studies. *Brief Bioinform* 2018, 19(4):700-712.

688 51. Yu J, Pressoir G, Briggs WH, Vroh Bi I, Yamasaki M, Doebley JF, McMullen MD, Gaut  
689 BS, Nielsen DM, Holland JB *et al*: A unified mixed-model method for association  
690 mapping that accounts for multiple levels of relatedness. *Nat Genet* 2006, 38(2):203-  
691 208.

692

693

694

695

696

697

698

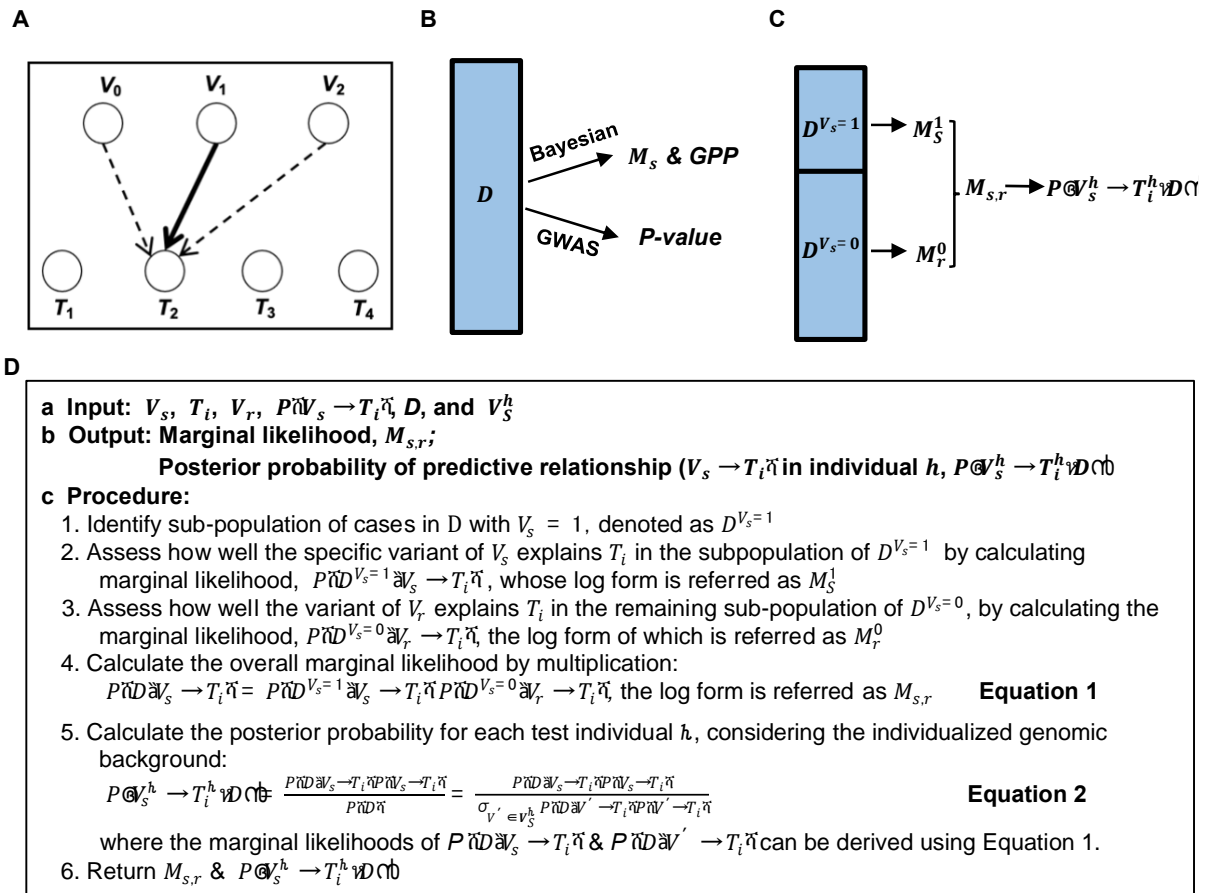
699

700

701

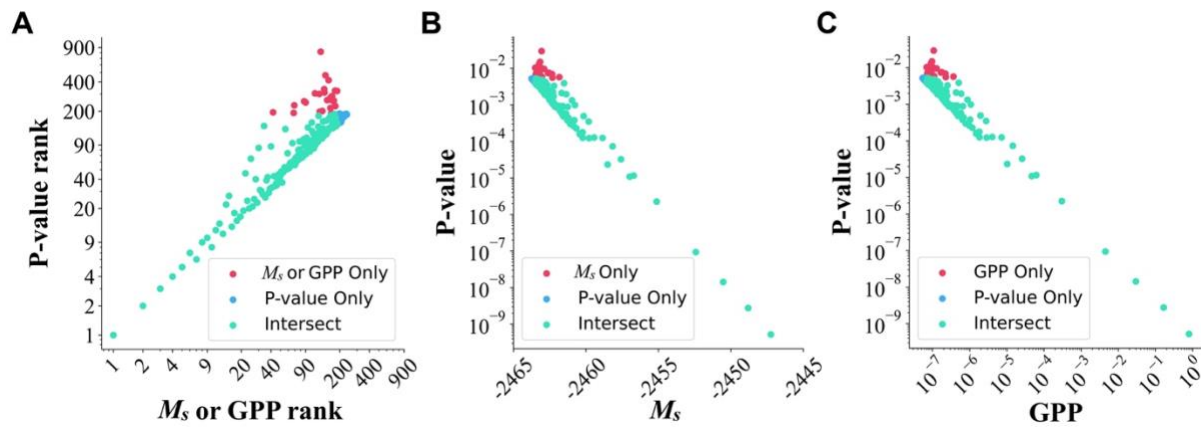
702

703 **Figures**



704

705 **Figure 1. The IBI Algorithm.** **A.** IBI uses a bipartite BN to model the probabilistic  
 706 relationships from genomic Variants to Traits.  $V$  nodes denote variants and  $T$  nodes denote  
 707 traits; the arcs denote the predictive relationships from  $V$  to  $T$  with only one assumed  
 708 predictive variant being evaluated at a time indicated by the solid arc. **B.** Using the entire  
 709 dataset ( $D$ ) or population to evaluate the association between a particular genomic variant  
 710 and a trait, GWAS methods output the p-value while the Bayesian method uses the marginal  
 711 likelihood ( $M_s$ ) and global posterior probability (GPP). **C.** Based on the value of a particular  
 712 variant  $V_s$ , IBI partitions the whole population into two subpopulations,  $D^{V_s=1}$  and  $D^{V_s=0}$ , and  
 713 derives the subpopulation-specific marginals,  $M_s^1$  and  $M_r^0$ , using  $V_s$  and  $V_r$  as the assumed  
 714 cause specifically. The overall marginal  $M_{s,r}$  and the individual-specific posterior probability,  
 715  $P(V_s^h \rightarrow T_i^h | D)$  for the SNP  $V_s$  can be further derived. **D.** Pseudo code for the IBI algorithm.



716

717 **Figure 2. A Bayesian method for GWAS analysis. A.** The p-value ranks and  $M_s$  or GPP

718 ranks are the same or similar for the top 189 SNPs selected by  $M_s$  or p-value ranking. **B.**

719 The p-values were very much negatively correlated with the  $M_s$  values of the top 189 SNPs

720 selected by  $M_s$  or p-value ranking. **C.** The p-values were very much negatively correlated

721 with the global posterior probabilities (GPP) of the top 189 SNPs selected by  $M_s$  or p-value

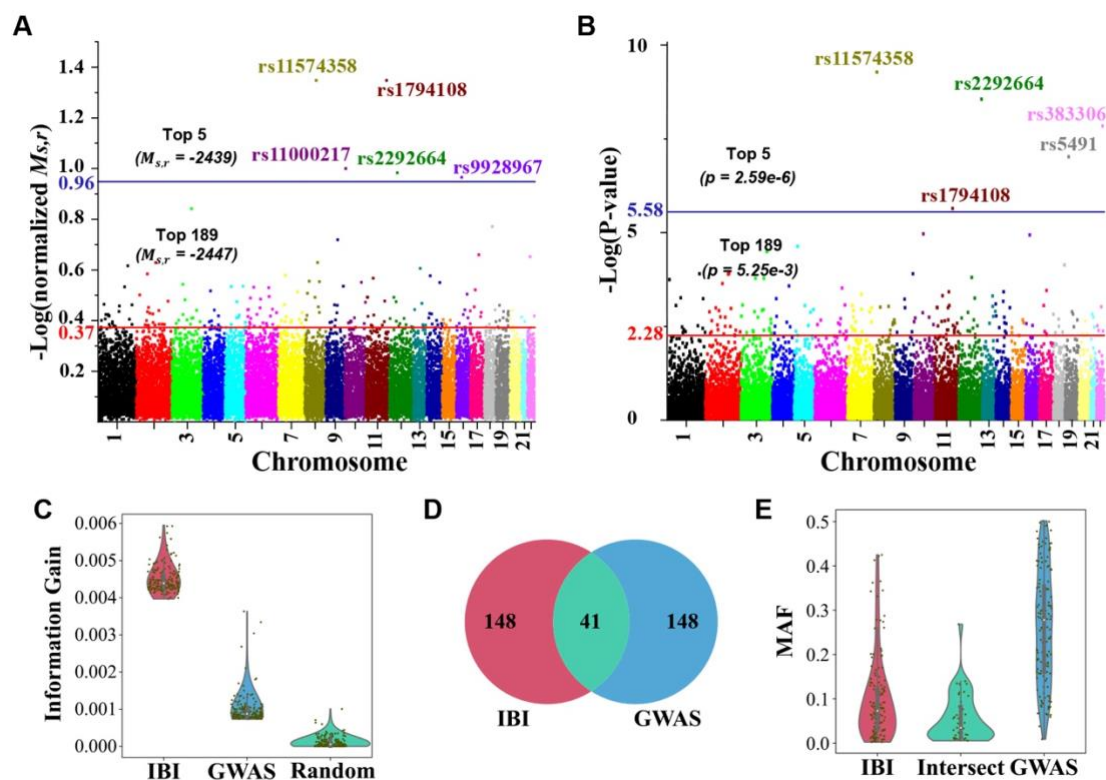
722 ranking.

723

724

725

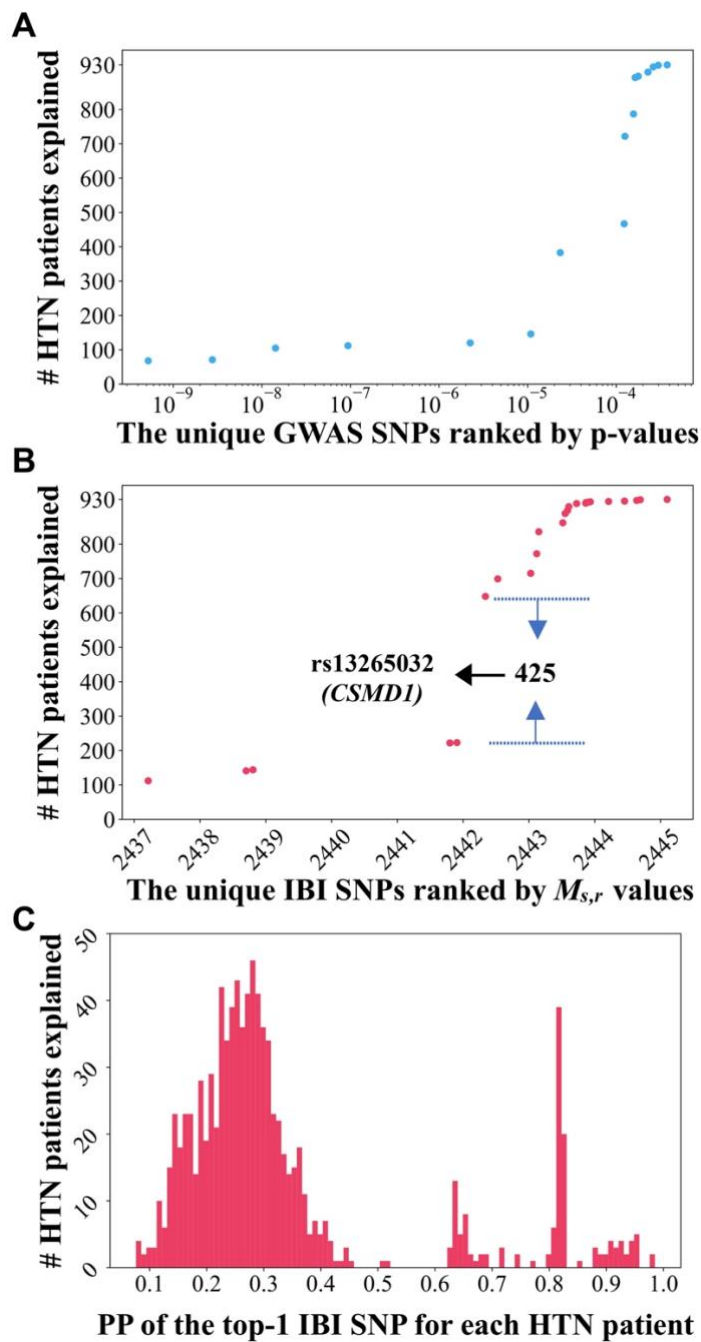
726



727

728 **Figure 3. Comparison of IBI and GWAS.** **A.** A Manhattan plot of SNPs' chromosome  
 729 location and normalized  $M_{s,r}$  values acquired by IBI. The threshold lines in blue and red  
 730 together with the corresponding threshold  $M_{s,r}$  values were labeled for both top 5 and top  
 731 189 SNPs ranked by  $M_{s,r}$ . Top five SNPs were annotated with rs IDs. **B.** A Manhattan plot of  
 732 SNPs' chromosome location and p-values acquired by GWAS. The threshold lines in blue  
 733 and red together with the corresponding threshold p-values were labeled for both top 5 and  
 734 top 189 SNPs ranked by p-values. Top five SNPs were annotated with rs IDs. **C.** Information  
 735 gain from top 189 SNPs ranked by IBI, GWAS and randomly-selected 189 SNPs. The black  
 736 dots represent the information gain values for individual SNPs. **D.** A Venn diagram of top  
 737 189 SNPs ranked by IBI and GWAS. **E.** Violin plots of the MAF distributions of the SNPs in  
 738 the three sections of D, IBI, Intersection, GWAS. The black dots represent the MAF values  
 739 for individual SNPs. The thick vertical gray bars show the interquartile range and the three  
 740 white dots represent the medians. Wider sections suggest higher probability for the given  
 741 MAF values.





742

743 **Figure 4. HTN patient coverage. A.** The 16 unique GWAS SNPs ranked by p-values were

744 plotted against the cumulative number of HTN patients explained. **B.** The 25 unique IBI

745 SNPs ranked by  $M_{s,r}$  were plotted against the cumulative number of HTN patients explained.

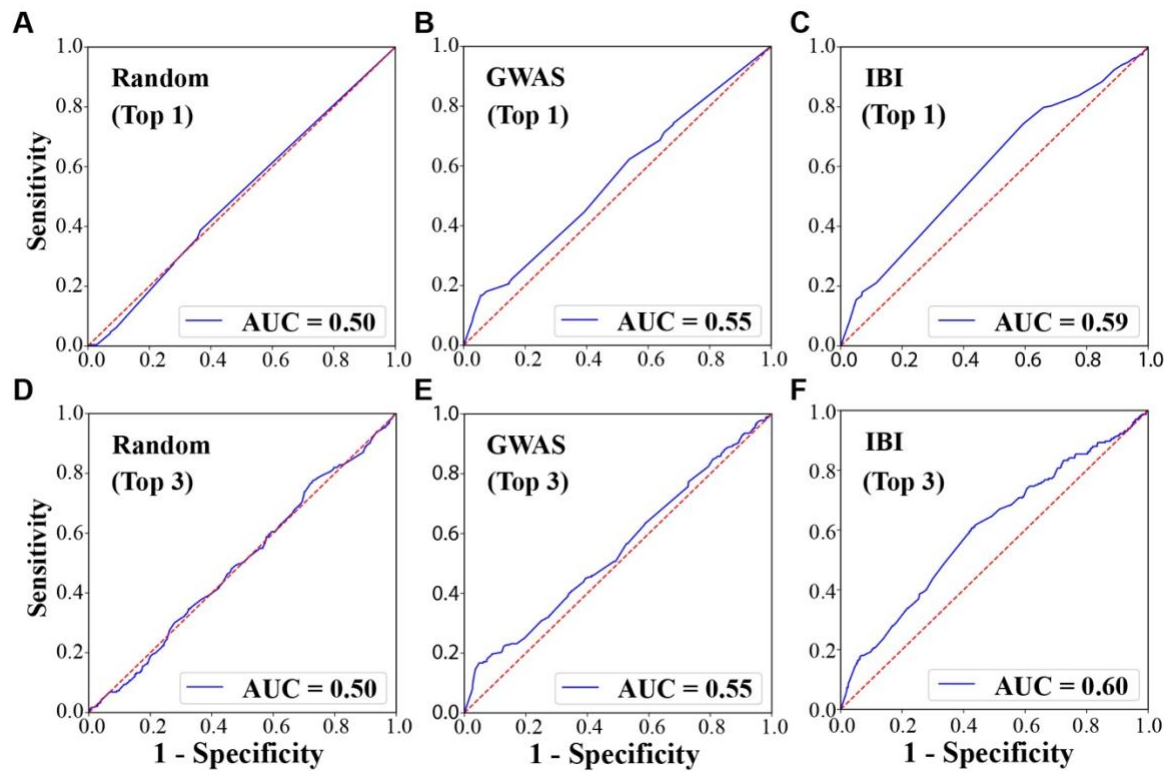
746 The number of HTN patients covered by each of these SNPs can be derived by taking the

747 difference of the two adjacent cumulative numbers of explained HTN patients as showed for

748 rs13265032 covering 425 HTN patients. **C.** Histogram of the individualized posterior

749 probability of the top-1 SNP assigned by IBI to each of the 930 HTN patient.





750

751 **Figure 5. ROC curves for predicting hypertension from top SNPs.** Top 1 SNP (A, B, C)  
752 and top 3 SNPs (D, E, F) selected by GWAS ( $p$ -value ranking), IBI ( $M_{S,r}$  ranking) or  
753 randomly were used to calculate the genetic risk scores for each test patient, which were  
754 further used to derive the AUROC for predicting hypertension. For these two experiments,  
755 random SNPs have AUROC as 0.50 (Top 1/3 SNP; (A, D)); GWAS-selected top SNPs have  
756 AUROC as 0.55 (B, E); IBI-selected top SNPs have AUROC as 0.59 for top 1 SNP and 0.60  
757 for top 3 SNPs.