

1 ***GSTM1* copy number is not associated with risk of kidney failure in**
2 **a large cohort**

3 Yanfei Zhang¹, PhD, Waleed Zafar², MD, Dustin N. Hartzel³, BS, Marc S. Williams¹, MD, Adrienne Tin⁴,
4 PhD, Alexander R. Chang^{2*}, MD, Regeneron Genetics Center⁵, Ming Ta M. Lee^{1*}, PhD on behalf of the
5 DiscovEHR collaboration

- 6 1. Genomic Medicine Institute, Geisinger, Danville, PA, USA
- 7 2. Kidney Institute, Geisinger, Danville, PA, USA
- 8 3. Phenomic Analytics & Clinical Data Core, Geisinger, Danville, PA, USA
- 9 4. Department of Epidemiology, Johns Hopkins University Bloomberg School of Public Health,
10 MD, USA
- 11 5. Regeneron Genetics Center, Regeneron Pharmaceuticals, Inc., Tarrytown, NY, USA

12 Drs. Zhang and Zafar contributed equally to this article, so do Drs. Chang and Lee.

13 * Corresponding: mlee2@geisinger.edu

14 **Corresponding author:** Ming Ta M. Lee, mlee2@geisinger.edu, Weis Center for Research, Geisinger,
15 100 North Academy Ave, Danville, PA 17822, USA. Fax: 570-214-7342. Tel: 570-271-6664.

16 The sources of support that require acknowledgment: The Regeneron Genetics Center funded the
17 collection of study samples, the generation of whole exome sequencing data.

18 **Running headline:** *GSTM1* and risk of kidney failure

19

20 **Abstract**

21 Deletion of *glutathione S-transferase μ1 (GSTM1)* is common in populations and has been asserted to
22 associate with chronic kidney disease progression in some research studies. The association needs to be
23 validated. We estimated *GSTM1* copy number using whole exome sequencing data in the DiscovEHR
24 cohort. Kidney failure was defined as requiring dialysis or receiving kidney transplant using data from the
25 electronic health record and linkage to the United States Renal Data System, or the most recent eGFR <
26 15 ml/min/1.73m². In a cohort of 46,983 unrelated participants, 28.8% of blacks and 52.1% of whites had
27 0 copies of *GSTM1*. Over a mean of 9.2 years follow-up, 645 kidney failure events were observed in
28 46,187 white participants, and 28 in 796 black participants. No significant association was observed
29 between *GSTM1* copy number and kidney failure in Cox regression adjusting for age, sex, BMI, smoking
30 status, genetic principal components, or co-morbid conditions (hypertension, diabetes, heart failure,
31 coronary artery disease, and stroke), whether using a genotypic, dominant, or recessive model. In
32 sensitivity analyses, *GSTM1* copy number was not associated with kidney failure in participants that were
33 45 years or older at baseline, had baseline eGFR < 60 ml/min per 1.73 m², or with baseline year between
34 1996-2002. In conclusion, we found no association between *GSTM1* copy number and kidney failure in a
35 large cohort study.

36 **Keywords:** *GSTM1*, kidney failure, copy number, large cohort

37

38 **Translational Statement**

39 Deletion of *GSTM1* has been shown to be associated with higher risk of kidney failure. However,
40 inconsistent results have been reported. We used electronic health record and whole exome sequencing
41 data of a large cohort from a single healthcare system to evaluate the association between *GSTM1* copy
42 number and risk of kidney failure. We found no significant association between *GSTM1* copy number and
43 risk of kidney failure overall, or in multiple sensitivity and subgroup analyses.

44

45

46 **Introduction**

47 *Glutathione S-transferase μ 1 (GSTM1)*, belongs to the family of glutathione-S-transferases that
48 metabolize a broad range of reactive oxygen species and aldehydes.^{1,2} Loss of *GSTM1* is very common in
49 the population; approximately 50% of whites and 27% of blacks have zero copies of *GSTM1*.³ Deletion of
50 one or both copies of the gene results in reduced amount of GSTM1,⁴ and could lead to increased
51 oxidative stress due to diminished ability to neutralize reactive chemical species. An association between
52 loss of *GSTM1* and chronic kidney disease (CKD) progression has been reported in the African American
53 Study of Kidney and Hypertension (AASK).⁵ The association between *GSTM1* copy number and kidney
54 failure was reported in the Atherosclerosis Risk in Communities (ARIC) study which had a larger sample
55 size, both in blacks and whites.⁶ However, data from smaller case-control studies have had mixed
56 findings.⁷⁻¹¹ The inconsistent results may be due to insufficient sample size and different populations and
57 disease background.

58 The Geisinger MyCode® Community Health Initiative is an EHR-linked biobank for precision medicine
59 research.¹² The population served by Geisinger have low rates of out migration, thus the electronic health
60 record (EHR) data are relatively complete longitudinally, with a median of 14 years of follow-up.¹² Also,
61 a high proportion of participants were found to have first and second-degree relatedness through cryptic
62 relatedness analysis.¹³ Through the ongoing DiscovEHR collaboration with the Regeneron Genetic
63 Center, whole exome sequence (WES) data are available from approximately 92,000 MyCode®
64 participants to date.^{12,14,15} These comprehensive clinical data are linked to matched genetic data and
65 provide power to identify disease-genetic associations.¹⁶

66 Confirming whether or not loss of *GSTM1* increases the risk of kidney failure is important, given the high
67 prevalence of *GSTM1* loss in the population and the serious morbidity and mortality associated with
68 kidney failure. In this study, we examine the relationship between *GSTM1* copy number and incident
69 kidney failure in the Geisinger MyCode cohort.

70 **Results**

71 **Study cohort and baseline characteristics**

72 There were 46,983 unrelated participants included in the main analysis. Among these, 46,187 (98.3%)
73 were white, and 796 (1.7%) were black. Frequency of *GSTM1* copy number differed between blacks and
74 whites (Supplemental Figure 4). Almost half of the blacks (49.2%) had 1 copy of *GSTM1* compared with
75 only 39.6% of whites, and more than half (52.1%) of whites had 0 copies of *GSTM1* compared with
76 28.8% of blacks. Chi-square test indicated *GSTM1* copy number follows the HWE (Table 1. P-values
77 were 1×10^{-4} for whites and 0.618 for blacks).

78 Baseline demographic and clinical characteristics of the participants stratified by race and *GSTM1* copy
79 numbers are provided in table 1. The average baseline age for whites and blacks was 51 and 44 years old
80 respectively. The mean baseline eGFR were 91 and 92 ml/min per 1.73 m^2 . Prevalence of hypertension
81 was 29% in whites and 44% in blacks, and prevalence of type 2 diabetes was 13% in whites and 21% in
82 blacks at baseline. No statistically significant differences were observed for baseline characteristics
83 among the three genotype groups when adjusted for multiple comparisons.

84 **Time-to-event analysis**

85 Over a mean follow-up of 9.2 years, there were a total of 645 kidney failure events in 46,187 white
86 participants. Over a mean follow-up of 5.5 years, there were 28 events in 796 blacks. No significant
87 difference in kidney failure-free survival was found across *GSTM1* copy numbers by log rank test (P=0.9
88 and 0.5 for whites and blacks, respectively, Figure 1).

89 Findings were similar in Cox regression analyses adjusting for other covariates (Table 2). In the genotypic
90 model, after adjusting for age, sex and the first 4 PCs, there was no difference in risk of kidney failure
91 among white participants with 0 copies (hazard ratio [HR] 1.01, 95% confidence interval [CI]: 0.81-1.48)
92 or 1 copy (HR 1.08, 95% CI: 0.80-1.47), compared to those with 2 copies. Findings were similar for black
93 participants in the genotypic model adjusting for age, sex, and the first 4 PCs (0 copies HR 0.68, 95% CI:
94 0.22-2.15; 1 copy HR 1.13, 95% CI: 0.44-2.90). In the dominant model, after adjusting for age, sex, and
95 the first 4 PCs, there was again no difference in risk of kidney failure among participants with 0 or 1

96 copies for white participants (HR 1.10, 95% CI: 0.82-1.46) or black participants with 0 or 1 copies (HR
97 0.96, 95% CI: 0.39-2.34), compared to their counterparts with 2 copies. In the recessive model, after
98 adjusting for age, sex, and the first 4 PCs, there was no difference in risk of kidney failure among white
99 participants with 0 copies (1.03, 95% CI: 0.88-1.20), or black participants with 0 copies (HR 0.63, 95%
100 CI: 0.25-1.56), compared to their counterparts with 1 or 2 copies. Adjustment for additional renal risk
101 factors in Model 2 yielded similar results (Table 2 and Supplemental Tables 1-3).

102 **Sensitivity analyses**

103 Analyses were performed on white participants who were older than 45 years of age at baseline
104 (n=31406), had baseline eGFR < 60 ml/min/1.73m² (n=3129), or had index date before 2003 (n=14572).
105 No statistically significant differences in the risk of kidney failure were identified for participants older
106 than 45 years of age at baseline (0 copies: HR 1.03, 95% CI: 0.74-1.42; 1 copy: HR 0.99, 95% CI: 0.71-
107 1.37), or those with baseline eGFR < 60 ml/min/1.73m² (0 copies: HR 1.13, 95% CI: 0.75-1.69; 1 copy
108 HR 1.04, 95% CI: 0.69-1.57), or for those with baseline year before 2003 (0 copies: HR 0.88, 95% CI:
109 0.59-1.32; 1 copy HR 0.88, 95% CI: 0.59-1.33) (Table 2 and Supplemental Tables 1-3).

110 **Discussion**

111 To our knowledge, this is the largest study to investigate the association of *GSTM1* copy number variation
112 with kidney failure. There were a total of 46,187 unrelated whites and 796 blacks in the study with an
113 average of 9.3 years' follow-up. The frequency of the *GSTM1* copy numbers were very similar to those
114 reported previously.^{3,5,6} Our results showed no significant association between *GSTM1* copy number and
115 risk of kidney failure in unadjusted or adjusted analyses whether using a genotypic, dominant, or
116 recessive genetic model.

117 Data from the ARIC and AASK cohorts suggested that the loss of *GSTM1* increased the risk of kidney
118 failure or accelerated CKD progression.^{5,6,17} In ARIC, a community-based cohort of middle-aged black
119 and white participants, there were 3461 white participants with WES reads. In fully adjusted models
120 among whites in ARIC, loss of *GSTM1* was associated with risk of kidney failure (0 or 1 copy vs. 2
121 copies: HR 2.54; 95% CI: 1.32-4.88). It is possible that differences in baseline characteristics of the study

122 populations could explain differing findings. Our cohort had a lower prevalence of smoking (never
123 smoking 48% vs. 38%), a higher prevalence of diabetes (14% vs. 8%), and was more contemporary
124 (median baseline year 2004 vs. ARIC baseline year 1987-1989). However, in sensitivity analyses in
125 14572 white participants with baseline year between 1996-2002 in our cohort, there was no association
126 between *GSTM1* loss and risk of kidney failure.

127 There is also the possibility that *GSTM1* may only be deleterious when GFR falls below certain levels, or
128 in the setting of specific types of kidney disease. In the AASK study, a randomized trial comparing
129 different antihypertensive medications and levels of blood pressure control in blacks with CKD attributed
130 to hypertension, 692 participants had *GSTM1* genotyping completed. Loss of *GSTM1* in AASK was
131 associated with increased risk of CKD progression (HR 0 copy: HR 1.88; 95% CI: 1.07-3.30; HR 1 copy:
132 1.68; 95% CI: 1.00-2.84). While there were few blacks in our study population with CKD, we found no
133 association between loss of *GSTM1* and risk of kidney failure in 3129 white participants with baseline
134 eGFR < 60 ml/min/1.73m². Some smaller case-control studies also tested a recessive genetic model (0
135 copies vs. 1 or 2 copies) of *GSTM1* and found that *GSTM1* 0 copy was associated with kidney failure in
136 several case-control studies.^{5,10,11,18} In our study, however, we examined both recessive and dominant
137 genetic models, and found no associations between *GSTM1* copy number and risk of kidney failure.

138 Some limitations of this study are worth noting. First, the genotypes of *GSTM1* in this study were derived
139 from WES and the thresholds were determined empirically from the histogram. Multiplex PCR validation
140 was not performed for this study. However, this method was used in the ARIC study, which showed
141 99.3% agreement when comparing the results of 0 copy with ≥ 1 copy from PCR.⁶ Moreover, the
142 frequency of *GSTM1* copy numbers in this study is similar to those reported previously for both whites
143 and blacks,^{5,6} and it follows the Hardy-Weinberg equilibrium. We also examined the possibility of miss-
144 mapping due to the complexity of *GSTM* locus. The histogram of the sum of normalized coverage of
145 *GSTM5*, the paralog gene of *GSTM1*, is unimodal distribution. Based on these evidences, we believe the
146 copy number results derived from WES had a low error rate and can be used with confidence for analysis.

147 Second, the sample size for blacks was small with only 796 patients available for analysis, limiting the

148 power to examine an association between *GSTM1* and kidney failure in this population. The strength of
149 this study is the large sample size of whites, which included 46,187 patients, allowing for large subgroup
150 analyses in older patients, those with baseline eGFR < 60 ml/min/1.73m², and those with longer duration
151 of follow-up. Finally, *GSTM1* deletion is known to be protective for lung cancer among individuals with a
152 high intake of cruciferous vegetable and deleterious among those with low intake of cruciferous
153 vegetable.¹⁹⁻²³ We were not able to analyze the effect modification between *GSTM1* copy number and
154 cruciferous vegetable intake.

155 In conclusion, our study does not support an association between loss of *GSTM1* and increased risk of
156 kidney failure. Additional research is needed to confirm whether loss of *GSTM1* increases risk of kidney
157 failure in certain subgroups such as blacks, and the interaction with diet.

158

159 **Methods**

160 **Study Population**

161 The study population included 72,756 participants in the Geisinger-Regeneron DiscovEHR cohort with
162 WES data and at least 3 serum creatinine values. We excluded participants who had baseline eGFR < 15
163 ml/min/1.73m² or history of dialysis or transplant (n=417), baseline age <18 (n=1,960), missing BMI
164 values (n=230), and unknown smoking status (n=2,694). Cryptic relatedness was evaluated using IBD
165 method in PLINK1.9.²⁴ One of the related pair of participants with PI_HAT >= 0.125 were removed to
166 reduce confounding from other shared genetic/environmental factors that could not be assessed, resulting
167 in 46983 participants (Supplement Figure 1). We estimated eGFR using the CKD-Epidemiology
168 Collaboration equation.^{25,26} We defined the study period, from a baseline time of the second serum
169 creatinine to the time of a renal endpoint, or the last serum creatinine test. All participants provided their
170 informed written consent, and the study was approved by the Geisinger Institutional Review Board.

171 **Outcome definition**

172 An assignment of kidney failure was made if any of the following criteria was met: 1) the last available
173 eGFR was less than 15 mL/min/1.73 m²; 2) the EHR showed an International Classification of Disease
174 (ICD) code for end stage renal disease (ESRD) (ICD9: 585.6, ICD10: N18.6); 3) receipt of dialysis or
175 transplant per linkage to the United States Renal Data System (USRDS). For participants who met more
176 than one criteria, the earliest documented date of the criterion was considered the kidney failure date.

177 **Clinical variables**

178 The following data elements were extracted from the EHR: baseline age, sex, self-reported race, body
179 mass index (BMI, BMI = Weight(kg)/Height(m)²), serum creatinine, smoking status; and ICD-9/10 coded
180 diagnosis of hypertension, diabetes, coronary artery disease, heart failure, and stroke. Genetic principal
181 components were calculated by PLINK1.9 (www.cog-genomics.org/plink/1.9/).²⁴

182 ***GSTM1* and *GSTM5* copy number estimation**

183 DNA was sequenced in two batches at the Regeneron Genetic Center. The WES data processing has been
184 described elsewhere,^{13,15} and details are also provided in the supplementary methods. Copy number of
185 *GSTM1* was estimated using the method reported previously.⁶ Briefly, sequence coverage was normalized
186 using CLAMMS software.²⁷ Normalized coverage of all eight exons of *GSTM1* was summed for each
187 participant. Thresholds for *GSTM1* copy number were determined empirically according to the
188 distribution of sum of normalized coverage (Supplementary Figure 2). The copy number of *GSTM5* was
189 estimated the same way as *GSTM1* (Supplementary notes).

190 **Statistical analyses**

191 Hardy-Weinberg equilibrium (HWE) was estimated using Chi-Square test for *GSTM1* copy number. A
192 more stringent significance level of 1e⁻⁵ was selected due to the large sample size (46,187) for Whites.
193 Nominal significance level of 0.05 was used for Blacks as the small sample size (796). Missing baseline
194 BMI values were imputed using mean BMI values calculated from all available BMI values for each
195 participant, and missing smoking status at baseline was imputed using the most recent recorded smoking
196 status in the EHR (Supplemental Figure 3). Baseline characteristics were compared using ANOVA for

197 continuous variables, and chi-squared tests for categorical variables. A p-value of $< 0.05/N$ was
198 considered statistically significant after Bonferroni adjustment for multiple comparisons, where N is the
199 number of variables compared. Kaplan–Meier curves were plotted by *GSTM1* copy number, and survival
200 differences between genotype groups were assessed using a log rank test. All subsequent analyses were
201 stratified by race given the difference in allele frequency and sample size. Two models were evaluated. In
202 Model 1, Cox proportional hazards model was adjusted for age, sex and the first four genetic principle
203 components; Based on Model 1, Model 2 was additionally adjusted for risk factors including baseline
204 eGFR, smoking status, BMI, hypertension, diabetes, coronary artery disease, heart failure and stroke.
205 Three genetic models were evaluated: 1) genotypic model, using copy number 0,1, and 2 as three
206 categories; 2) dominant model (0 or 1 copy vs. 2 copies *GSTM1*); 3) recessive model (1 or 2 copies vs. 0
207 copies of *GSTM1*). To explore whether the effect of *GSTM1* loss was stronger in specific higher-risk
208 subgroups, sensitivity analyses were completed including subsets of 1) older participants (baseline age \geq
209 45 years); 2) participants with CKD (baseline eGFR < 60 ml/min per 1.73 m²); 3) participants with longer
210 follow-up (baseline year 1996-2002). Power was estimated using powerCT() function in R package
211 powerSurvEpi 0.1.0. Given the current sample size and event number, we have power of ≥ 0.8 to test
212 hazard ratio of at least 1.4 and 2.8 for white and black cohort, respectively. All analyses were performed
213 using R (version 3.4.3).

214 **Disclosure**

215 Regeneron Genetics Center authors work for Regeneron Pharmaceuticals.

216 **Acknowledgments**

217 The authors thank the staff and participants of MyCode. The Regeneron Genetics Center funded the
218 collection of study samples, the generation of whole exome sequencing data. Geisinger provided funding
219 for clinical data extraction and other analysis. A.C. is supported by National Institutes of Health/National
220 Institute of Diabetes and Digestive and Kidney Diseases grant K23 DK106515-01. Data reported here
221 were supplied by the U.S. Renal Data System. The interpretation and reporting of these data are the

222 responsibility of the authors and in no way should be seen as an official policy or interpretation of the
223 U.S. government. The authors thank Ilene G. Ladd who is the project manager at Genomic Medicine
224 Institute for English editing.

225

226 **References**

- 227 1. Yang Y, Parsons KK, Chi L, Malakauskas SM, Le TH. Glutathione S-transferase-micro1 regulates
228 vascular smooth muscle cell proliferation, migration, and oxidative stress. *Hypertension*.
229 2009;54(6):1360-1368.
- 230 2. Hayes JD, Flanagan JU, Jowsey IR. Glutathione transferases. *Annu Rev Pharmacol Toxicol*.
231 2005;45:51-88.
- 232 3. Garte S, Gaspari L, Alexandrie AK, et al. Metabolic gene polymorphism frequencies in control
233 populations. *Cancer Epidemiol Biomarkers Prev*. 2001;10(12):1239-1248.
- 234 4. Board P, Coggan M, Johnston P, Ross V, Suzuki T, Webb G. Genetic heterogeneity of the human
235 glutathione transferases: a complex of gene families. *Pharmacol Ther*. 1990;48(3):357-369.
- 236 5. Chang J, Ma JZ, Zeng Q, et al. Loss of GSTM1, a NRF2 target, is associated with accelerated
237 progression of hypertensive kidney disease in the African American Study of Kidney Disease
238 (AASK). *Am J Physiol Renal Physiol*. 2013;304(4):F348-355.
- 239 6. Tin A, Scharpf R, Estrella MM, et al. The Loss of GSTM1 Associates with Kidney Failure and Heart
240 Failure. *J Am Soc Nephrol*. 2017;28(11):3345-3352.
- 241 7. Yang Y, Kao MT, Chang CC, et al. Glutathione S-transferase T1 deletion is a risk factor for
242 developing end-stage renal disease in diabetic patients. *Int J Mol Med*. 2004;14(5):855-859.
- 243 8. Tiwari AK, Prasad P, B KT, et al. Oxidative stress pathway genes and chronic renal insufficiency in
244 Asian Indians with Type 2 diabetes. *J Diabetes Complications*. 2009;23(2):102-111.
- 245 9. Nomani H, Hagh-Nazari L, Aidy A, et al. Association between GSTM1, GSTT1, and GSTP1 variants
246 and the risk of end stage renal disease. *Ren Fail*. 2016;38(9):1455-1461.
- 247 10. Gutierrez-Amavizca BE, Orozco-Castellanos R, Ortiz-Orozco R, et al. Contribution of GSTM1,
248 GSTT1, and MTHFR polymorphisms to end-stage renal disease of unknown etiology in Mexicans.
249 *Indian J Nephrol*. 2013;23(6):438-443.
- 250 11. Agrawal S, Tripathi G, Khan F, Sharma R, Baburaj VP. Relationship between GSTs gene
251 polymorphism and susceptibility to end stage renal disease among North Indians. *Ren Fail*.
252 2007;29(8):947-953.
- 253 12. Carey DJ, Fetterolf SN, Davis FD, et al. The Geisinger MyCode community health initiative: an
254 electronic health record-linked biobank for precision medicine research. *Genet Med*.
255 2016;18(9):906-913.
- 256 13. Staples J, Maxwell EK, Gosalia N, et al. Profiling and Leveraging Relatedness in a Precision
257 Medicine Cohort of 92,455 Exomes. *Am J Hum Genet*. 2018;102(5):874-889.
- 258 14. Abul-Husn NS, Manickam K, Jones LK, et al. Genetic identification of familial
259 hypercholesterolemia within a single U.S. health care system. *Science*. 2016;354(6319).
- 260 15. Dewey FE, Murray MF, Overton JD, et al. Distribution and clinical impact of functional variants in
261 50,726 whole-exome sequences from the DiscovEHR study. *Science*. 2016;354(6319).
- 262 16. Rader DJ, Damrauer SM. "Pheno"menal value for human health. *Science*. 2016;354(6319):1534-
263 1536.

- 264 17. Bodonyi-Kovacs G, Ma JZ, Chang J, et al. Combined Effects of GSTM1 Null Allele and APOL1 Renal
265 Risk Alleles in CKD Progression in the African American Study of Kidney Disease and
266 Hypertension Trial. *J Am Soc Nephrol*. 2016;27(10):3140-3152.
- 267 18. Suvakov S, Damjanovic T, Stefanovic A, et al. Glutathione S-transferase A1, M1, P1 and T1 null or
268 low-activity genotypes are associated with enhanced oxidative damage among haemodialysis
269 patients. *Nephrol Dial Transplant*. 2013;28(1):202-212.
- 270 19. London SJ, Yuan JM, Chung FL, et al. Isothiocyanates, glutathione S-transferase M1 and T1
271 polymorphisms, and lung-cancer risk: a prospective study of men in Shanghai, China. *Lancet*.
272 2000;356(9231):724-729.
- 273 20. Zhao B, Seow A, Lee EJ, et al. Dietary isothiocyanates, glutathione S-transferase -M1, -T1
274 polymorphisms and lung cancer risk among Chinese women in Singapore. *Cancer Epidemiol
275 Biomarkers Prev*. 2001;10(10):1063-1067.
- 276 21. Wang LI, Giovannucci EL, Hunter D, Neuberger D, Su L, Christiani DC. Dietary intake of Cruciferous
277 vegetables, Glutathione S-transferase (GST) polymorphisms and lung cancer risk in a Caucasian
278 population. *Cancer Causes Control*. 2004;15(10):977-985.
- 279 22. Brennan P, Hsu CC, Moullan N, et al. Effect of cruciferous vegetables on lung cancer in patients
280 stratified by genetic status: a mendelian randomisation approach. *Lancet*. 2005;366(9496):1558-
281 1560.
- 282 23. Carpenter CL, Yu MC, London SJ. Dietary isothiocyanates, glutathione S-transferase M1 (GSTM1),
283 and lung cancer risk in African Americans and Caucasians from Los Angeles County, California.
284 *Nutr Cancer*. 2009;61(4):492-499.
- 285 24. Chang CC, Chow CC, Tellier LC, Vattikuti S, Purcell SM, Lee JJ. Second-generation PLINK: rising to
286 the challenge of larger and richer datasets. *Gigascience*. 2015;4:7.
- 287 25. Levey AS, Stevens LA, Schmid CH, et al. A new equation to estimate glomerular filtration rate.
288 *Ann Intern Med*. 2009;150(9):604-612.
- 289 26. Levey AS, Stevens LA. Estimating GFR using the CKD Epidemiology Collaboration (CKD-EPI)
290 creatinine equation: more accurate GFR estimates, lower CKD prevalence estimates, and better
291 risk predictions. *Am J Kidney Dis*. 2010;55(4):622-627.
- 292 27. Packer JS, Maxwell EK, O'Dushlaine C, et al. CLAMMS: a scalable algorithm for calling common
293 and rare copy number variants from exome sequencing data. *Bioinformatics*. 2016;32(1):133-
294 135.

295

296

297

298

299 **Author Contributions**

300 A.T, A.C, and M.T.M.L designed the study; D.N.H extracted data; Y.Z and W.Z performed the analysis
 301 and drafted the manuscript. M.S.W, M.T.M.L and A.C did critical review on the paper. RGC supported
 302 the study and did review the manuscript. All authors approved the final version of the manuscript.

303 **Tables:**

304 **Table 1: Baseline characteristics of participants by *GSTM1* copy number**
 305

Characteristic	Whites				Blacks				
	Copy number	0	1	2	P	0	1	2	P
N.		24141	18255	3791	1E-4 ^	233	388	175	0.618 ^
Age		50.9(14.1)	50.8(14.0)	51.2(14.3)	0.317	41.8(12.3)	44.7(12.1)	44.3(13.2)	0.015
Female, n(%)		14184(58.8)	10750(58.9)	2208(58.2)	0.763	148(63.5)	240(61.9)	103(58.9)	0.629
eGFR		91.8(20.3)	91.7(20.4)	91.3(20.8)	0.326	95.0(22.9)	90.8(23.2)	91.9(24.4)	0.095
BMI		31.6(7.9)	31.6(7.8)	31.6(7.9)	0.817	33.5(9.3)	34.7(9.6)	34.5(10.2)	0.299
Smoking Status									
Current Smoker		4688 (19.4)	3521 (19.3)	769 (20.3)	0.418	65 (27.9)	109 (28.1)	55 (31.4)	0.785
Former Smoker		7803 (32.3)	5829 (31.9)	1231 (32.5)		64 (27.5)	117 (30.2)	52 (29.7)	
Never Smoker		11650 (48.3)	8905 (48.8)	1791 (47.2)		104 (44.6)	162 (41.8)	68 (38.9)	
Hypertension		7193 (29.8)	5409 (29.6)	1118 (28.4)	0.893	96 (41.2)	169 (43.6)	83 (47.4)	0.453
T1DM		251 (1.0)	205 (1.1)	50 (1.3)	0.277	6 (2.5)	12 (3.1)	5 (2.9)	0.932
T2DM		3130 (13.0)	2299 (12.6)	483 (12.7)	0.522	40 (17.2)	92 (23.7)	35 (20.0)	0.143
CAD		1384 (5.7)	1072 (5.9)	241 (6.4)	0.304	0	0	0	/
HF		568 (2.4)	402 (2.2)	110 (2.9)	0.034	10 (4.3)	12 (3.1)	7 (4.0)	0.713
Stroke		413(1.7)	304(1.7)	69(1.8)	0.789	4(1.7)	11(2.8)	8(4.6)	0.233

306 ^ Hardy-Weinberg equilibrium test using Chi-Square test. Consider the sample size for Whites (46,187), we selected a more stringent
 307 significance level of 1E-5.

308 P<0.0045 (0.05/11) was considered statistically significant. T1DM: Type 1 diabetes mellitus. T2DM: type 2 diabetes mellitus. CAD: coronary
 309 artery disease; HF: heart failure.

310

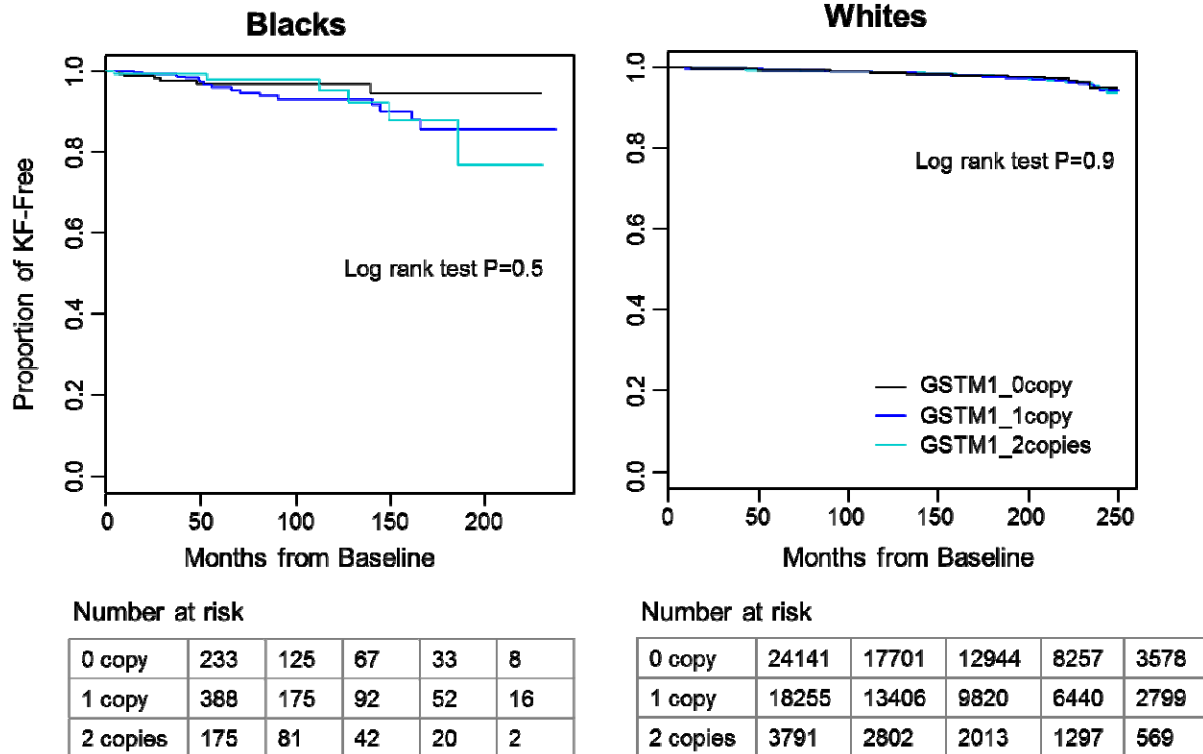
311 **Table 2: Risk of kidney failure associated with *GSTM1* copy number**

Population	Hazard Ratio [95% CI]							
	Genotypic Model (Ref: 2 copies)				Dominant Model		Recessive Model	
	0 copy	P	1 copy	P	0/1 vs 2	P	0 vs 1/2	P
Blacks (28/796)								
Model 1	0.68 [0.22, 2.15]	0.516	1.13 [0.44, 2.90]	0.802	0.96 [0.39, 2.34]	0.937	0.63 [0.25, 1.56]	0.316
Model 2	0.67 [0.19, 2.34]	0.527	0.63 [0.22, 1.84]	0.399	0.64 [0.23, 1.77]	0.391	0.93 [0.33, 2.59]	0.883
Whites (645/46187)								
Model 1	1.01 [0.81, 1.48]	0.537	1.08 [0.80, 1.47]	0.604	1.10 [0.82, 1.46]	0.551	1.03 [0.88, 1.20]	0.732
Model 2	1.14 [0.84, 1.54]	0.398	1.15 [0.84, 1.56]	0.384	1.14 [0.85, 1.53]	0.375	1.02 [0.87, 1.19]	0.836
Whites, Baseline Age>=45 (518/31406)								
Model 1	1.03 [0.74, 1.42]	0.878	0.99 [0.71, 1.37]	0.931	1.01 [0.74, 1.34]	0.960	1.04 [0.87, 1.23]	0.673
Model 2	1.12 [0.81, 1.54]	0.504	1.07 [0.77, 1.49]	0.677	1.10 [0.80, 1.50]	0.562	1.05 [0.89, 1.25]	0.553
Whites, Baseline eGFR<60 (304/3129)								
Model 1	1.13 [0.75, 1.69]	0.556	1.04 [0.69, 1.57]	0.860	1.09 [0.74, 1.61]	0.670	1.10 [0.87, 1.37]	0.429
Model 2	1.03 [0.68, 1.55]	0.891	1.03 [0.68, 1.56]	0.897	1.03 [0.69, 1.53]	0.889	1.01 [0.80, 1.27]	0.956
Whites, Baseline Year before 2003 (310/14572)								
Model 1	0.88 [0.59, 1.32]	0.542	0.88 [0.59, 1.33]	0.553	0.88 [0.60, 1.30]	0.529	0.98 [0.78, 1.22]	0.841
Model 2	0.86 [0.57, 1.29]	0.460	0.89 [0.59, 1.34]	0.563	0.87 [0.59, 1.29]	0.485	0.95 [0.76, 1.19]	0.648

312 Model 1: adjusted for age, sex, first 4 genetic principle components. Model 2: Model 1 plus baseline eGFR, BMI, smoking status, hypertension,
 313 type 1 diabetes, type 2 diabetes, heart failure, coronary artery disease, and stroke.

314

315 **Figures and figure legends:**
 316



317

318 **Figure 1:** Kaplan-Meier survival curves of time to kidney failure by *GSTM1* copy number in white and
 319 black participants. In both races, participants without *GSTM1* (black line) had longer kidney failure-free
 320 survival time than those with 1 or 2 copies of *GSTM1* (blue and cyan lines). No significant difference
 321 was found by log rank test. Tables below the figure showed the number of patients at risk with start time
 322 at Month 0, 50, 100, 150 and 200.

323

324