1    **Intrinsic DNA topology as a prioritization metric in genomic fine-mapping studies.**

2    Hannah C. Ainsworth,[1]* Timothy D. Howard,[2] and Carl D. Langefeld[1]**
3    [1]Department of Biostatistics and Data Science, Wake Forest School of Medicine, Winston-Salem, NC,
4    27157, USA.
5    [2]Department of Biochemistry, Wake Forest School of Medicine, Winston-Salem, NC, 27157, USA.
6     *Correspondence: hainswor@wakehealth.edu
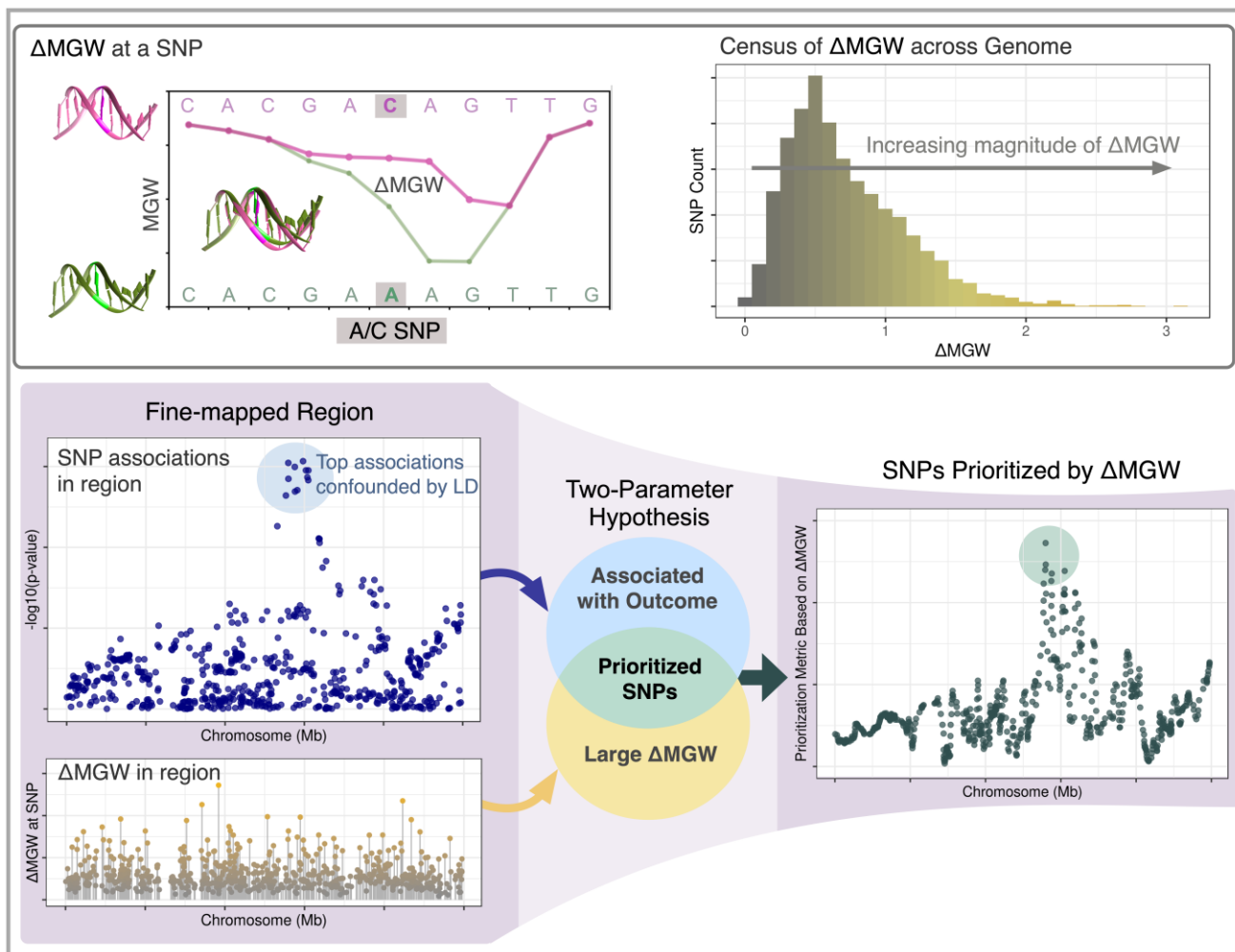7    **Correspondence: clangefe@wakehealth.edu
8
9
10

11   **Graphical Abstract**

12          We hypothesize that SNPs imposing dissimilar minor groove width profiles (ΔMGW) are

13   more likely to alter function.  ΔMGW was interrogated genome-wide and then used as a

14   weighting metric for fine-mapping associations.

15

16    **Abstract**

17        In genomic fine-mapping studies, some approaches leverage annotation data to

18    prioritize likely functional polymorphisms. However, existing annotation sources often present

19    challenges as many: lack data for novel variants, offer no context for noncoding regions, and/or

20    are confounded with linkage disequilibrium. We propose a novel annotation source – sequence-

21    dependent DNA topology – as a prioritization metric for fine-mapping. DNA topology and

22    function are well-intertwined, and as an intrinsic DNA property, it is readily applicable to any

23    genomic region. Here, we constructed and applied, Minor Groove Width (MGW), as a

24    prioritization metric. Using an established MGW-prediction method, we generated an MGW

25    census for 199,038,197 SNPs across the human genome. Summarizing a SNP's change in

26    MGW ($\Delta$MGW) as a Euclidean distance, $\Delta$MGW exhibited a strongly right-skewed distribution,

27    highlighting the infrequency of SNPs that generate dissimilar shape profiles. We hypothesized

28    that phenotypically-associated SNPs can be prioritized by $\Delta$MGW. We applied Bayesian and

29    frequentist MGW-prioritization approaches to three non-coding regions associated with System

30    Lupus Erythematosus in multiple ancestries. In two regions, including $\Delta$MGW resolved the

31    association to a single, trans-ancestral, SNP, corroborated by external functional data.

32    Together, this study presents the first usage of sequence-dependent DNA topology as a

33    prioritization metric in genomic association studies.

34

**Introduction**

Genetic association studies have successfully identified thousands of loci associated with a broad range of phenotypes.(1) However, despite the abundance of these genomic associations, analytic challenges have largely hindered identification of the specific genomic drivers of disease.(2–4) First, linkage disequilibrium (LD) constitutes a major analytic challenge, as highly correlated variants exhibit comparable evidence of association, making it difficult to statistically isolate causal polymorphisms. Second, many associated single nucleotide polymorphisms (SNPs) reside in non-coding regions, occluding functional relevance without additional context and information. Even with increased sample sizes and variant coverage, these challenges remain.(2–5) In-depth functional analyses are not practical for a large number of variants, and thus, there remains the need to effectively prioritize the most likely causal variants for follow-up studies and approaches (e.g. CRISPR).

To prioritize potential causal variants, association results can be weighted by external functional information (e.g. histone modifications, eQTL status, transcription factor binding sites).(5–8) This approach has been successful in reducing and refining associated variants, and there are a growing number of tools and methods that integrate external data with genomic association studies.(6, 9–13)  Nevertheless, such methods are not without limitations. Importantly, the choice of annotation and database bias are strong factors for consideration as missing or incomplete functional data could result in down-weighting potentially causal polymorphisms. These challenges particularly arise for regions with no (presently) known functional implications. Additionally, many annotation resources are based on European data; and thus may offer limited information for genetic studies in non-European individuals (e.g. novel regions).(14, 15) Such limitations can reduce the rate of progress in understanding the functional impact of ancestry-specific associations and perpetuate health disparities.(16, 17) To alleviate some of these biases imposed by external datasets, we propose a prioritization

60    approach that leverages information intrinsic to the DNA itself, sequence-dependent DNA

61    topology.

62        From chromatin conformation to selective protein binding,(18–26) DNA is a highly

63    dynamic macromolecule with structure inherently linked to function. Sequence-dependent DNA

64    topology (or shape) refers to the geometric parameters (measured in Angstroms or degrees)

65    between successive nucleotides in a DNA sequence.(24, 27–29) The sequence dependency of

66    these spatial measures (**Figure 1**) has been well-studied and in recent years, increasingly

67    connected to various functional implications, including protein binding, DNA stability, and

68    methylation.(18, 20, 21, 23, 30–38)  High-throughput DNA shape prediction methods now

69    enable exploration of DNA topology on a genome-wide scale, and thus, provide new

70    opportunities in association studies.(24, 39)

71        This study presents using sequence-dependent DNA topology as a prioritization metric

72    in genomic association studies. Here, we focused on minor groove width (MGW), which

73    measures the distance (Angstroms, Å) between the sugar phosphate backbone of the forward

74    and reverse strands. For each SNP, we analyzed its change in minor groove width (ΔMGW) to

75    evaluate whether the SNP's alleles created similar or divergent MGW profiles.  MGW has been

76    implicated in numerous protein binding studies and used in transcription factor binding

77    prediction algorithms.(18, 20, 24, 32, 34, 36, 37, 40, 41) Recently it was studied in the context of

78    purifying selection, where "shape disrupting variants" (examples shown in **Figures 2** and **3**) tend

79    to be less common in functional regions (shape-preserving polymorphisms being more

80    frequent).(42) Thus, we proposed that if a phenotypically-associated SNP also yields a large

81    ΔMGW, it is more likely to be causal as a function of divergent shape profiles.

82        We specifically hypothesized that highly correlated SNPs in a phenotype-associated

83    region can be functionally prioritized using each SNP's magnitude of ΔMGW.  We evaluated this

84    hypothesis in three stages. First, using an established MGW-prediction algorithm(39), we

85    generated the complete sample space for ΔMGW for all possible input sequences. Second, we

86    evaluated the observed frequency of ΔMGW across the human genome using bi-allelic SNPs in

87    the dbSNP SNP150 dataset. Third, we tested this approach by prioritizing SNPs in three

88    genomic regions previously associated with systemic lupus erythematosus (SLE)(43) leveraging

89    both frequentist and Bayesian association methods.

90    **Methods and Materials**

91

92    **Calculation of ∆MGW for a bi-allelic SNP.**

93        The predicted MGW for a given sequence was obtained using the DNAshapeR package

94    (https://bioconductor.org/packages/release/bioc/html/DNAshapeR.html) , available through

95    Bioconductor.(39) DNAshapeR calculates DNA features using Monte Carlo simulations for

96    nucleotide structure based on DNA sequence fragments. DNA feature predictions are based on

97    a rolling window of five nucleotides for a given n-length sequence. For this study, to capture the

98    MGW at a SNP, we used the four flanking (up and downstream) nucleotides (9-mer sequence)

99    as input. Each bi-allelic SNP produces two unique 9-mer sequences (one sequence for each

100   allele) and thus, both of a SNP's sequences were submitted to DNAshapeR to obtain the

101   corresponding feature vectors for MGW. The MGW was retained for the nucleotide at the SNP's

102   position as well as +/- 1 nucleotides. Capturing MGW for additional bases would require longer

103   input sequences, which could introduce additional variability (e.g. SNPs within the flanking

104   sequence). The ∆MGW was calculated as a Euclidean distance for the SNP and +/- 1 base

105   (**Figure 2**).

106   **Generation of ∆MGW sample space**

107       To calculate the entire sample space for ∆MGW, we generated a dataset of all possible

108   input sequences. Since our goal was to evaluate the ∆MGW at a SNP with +/- 4 base pairs,

109   input sequences required nine nucleotides.  Thus, all possible combinations of Adenine,

110   Cytosine, Guanine, and Thymine, generated 262,144 9-mer sequences. From this dataset, all

5

111    possible bi-allelic pairings (A/C, A/G, A/T, C/G, C/T, G/T) were created on the 5th nucleotide of

112    each sequence ("SNP position") while holding the flanking nucleotides constant, generating

113    393,216 9-mer pairings. These 9-mer pairings represent every possible sequence combination

114    that could be observed for a bi-allelic SNP (**Figure 3**). These paired sequences were evaluated

115    for ∆MGW using the previously described method.

116    **Visualization of DNA sequences**

117        DNA shape measures, provided by DNAshapeR, were submitted as a parameter file to

118    the 3D-Dart webportal (http://milou.science.uu.nl/services/3DDART/) for a 'BDNA nucleic acid'.

119    (44) Resulting pdb files from 3D-Dart were then visualized using Chimera

120    (https://www.cgl.ucsf.edu/chimera/).(45)

121    **Curating dbSNPs150 database**

122        The NCBI hg19 dbSNPs150 data file (snp150.txt.gz) was downloaded via UCSC

123    GoldenPath (hgdownload.cse.ucsc.edu) on July 6, 2018.(46) Insertion-deletions, tri-allelic,

124    quad- allelic, and multiple nucleotide polymorphisms were excluded. Retained bi-allelic SNPs

125    were limited to those located on chromosomes 1-22 and X. Any SNPs that were labeled with

126    "Unusual Conditions" as defined by UCSC were excluded, as these indicate possible

127    discrepancies among alleles and/or potential mapping issues (e.g. SNP flanking sequence

128    aligns to more than one location in the reference assembly).(46, 47) The pruned bi-allelic

129    dataset contained 199,038,272 SNPs.

130        For dbSNP 150 data, each SNP's flanking sequence of four nucleotides was retrieved

131    from the Human Reference Genome (downloaded October 2017)(48) using SAMTOOLS. For

132    each SNP, the dbSNP "Strand" variable was used to inform if the alleles reported by dbSNP

133    aligned with the reference genome. All SNPs were successfully queried against the reference

134    genome. There were 75 SNPs that contained at least one flanking base encoded as "N" (any

6

135  base) and were excluded from summarizations, leaving a final dataset of 199,038,197 SNPs.

136  The ΔMGW for these sequences were obtained as described above.

137
138  **SLE Immunochip Data for fine-mapping analyses**

139  Genomic data for fine-mapping analyses came from the published trans-ancestral SLE

140  Immunochip study; genotype calling and genomic quality control methods were previously

141  described.(43) This data includes three ancestries, European Ancestry (EA), African Ancestry

142  (AA), and Hispanic Ancestry (HA), with large case-control counts: EA (6,748; 11,516), AA

143  (2,970; 2,452), and HA (1,872; 2,016).

144  Genomic regions were named for the genes in physical proximity to the region of

145  association. Non-HLA genomic regions were selected for fine-mapping if the region contained

146  SNPs reaching genome-wide significance (p< $5\times10^{-8}$) in at least two ancestry-specific

147  analyses.(43) We also limited our analyses to regions where the top associations mapped to

148  non-coding regions (e.g. introns, intergeneic), where we hypothesize DNA topology might

149  provide novel insight to the fine-mapping analyses. Genomic regions containing *FAM167A-BLK*

150  (8p23), *STAT4* (2q32), and *TNIP1* (5q33) met these criteria. Quality controlled genomic data for

151  these regions were extracted using a 250 kb window around the previously reported top

152  association from the Immunochip analysis.(43)

153  SNPs from the selected genomic regions were queried against the human reference

154  genome to retrieve the four flanking bases. Each SNP's strand information (based on Illumina

155  Infinium Immunochip documentation) was utilized to ensure that the corresponding alleles

156  appropriately aligned with the reference genome.

157  **Statistical Analyses.**

158  *Single-SNP associations*. Single-SNP associations were previously reported and

159  described in the transancestral SLE Immunochip study.(43)

160    *SKAT analyses*. The previous single-SNP logistic regression analyses (43) did not

161    incorporate SNP-specific weights/information. Thus, SNPs in high LD yielded comparable

162    association values. The Sequence Kernel Association Test (SKAT) is a regression approach

163    that was designed to handle covariates and SNP-specific weights through a weighted linear

164    kernel.(49) It was shown that well-selected SNP weights can yield better statistical power (e.g.

165    increasing weight of functional variants).(49)  SKAT was originally developed to leverage minor

166    allele frequency (MAF), as the weighting scheme in rare variant studies; however, the SKAT

167    framework is a general method that can accommodate any user-specified SNP weights.(49)

168    Here, we used ∆MGW as the weighting scheme. A variation of SKAT is the Optimal unified test

169    which combines both SKAT and the burden test (SKAT-O).(12)  The SKAT-O test statistic is a

170    weighted average of the SKAT and burden test statistics and can be beneficial when applying to

171    genomic regions where one test may be better powered than another.(50)  Primary advantages

172    of burden tests occur when a large number of variants are causal and for smaller sample sizes

173    (SKAT loses power in small sample sizes, <2000 cases and controls). Generally, burden tests

174    do not perform as well as SKAT when a large proportion of the variants are non-causal.(12, 49,

175    50) In this study, our datasets are large (AA: 5,422; EA: 18,264; HA: 2,016), and we expect

176    many of the highly associated SNPs in LD to be non-causal; thus, in this scenario we selected

177    SKAT to be more appropriate, which is consistent with published power calculations and

178    simulations.(12, 49, 50)  SKAT was applied to genomic regions through its implementation in

179    the R package, SKAT (https://CRAN.R-project.org/package=SKAT). For each genomic region,

180    the model parameters and residuals were calculated for SKAT using SKAT_Null_Model() for a

181    dichotomous outcome (case/controls status) and previously described (43) population-specific

182    factors (to account for admixture). Since all datasets (AA, EA, and HA) had a sample size

183    greater than 2,000 cases and controls, no small-sample adjustment was applied. Within each

184    genomic region, adjacent 5-SNP windows were generated, offset by 1 SNP. Each window was

185    evaluated using the SKATbinary() with method=SKAT and a linear-weighted kernel with SNPs

8

186    weighted by their ∆MGW. To evaluate consistency of the results (e.g. for SNPs outside of the

187    main peak of association), genomic regions were also evaluated using equal-weighting for all

188    SNPs. Given the small window size (n=5 SNPs), we expect a large proportion of each window

189    to contain non-causal SNPs, further supporting our selection of SKAT. For comparison, we also

190    applied SKAT-O but noted minimal differences on the final outcome. To localize the top

191    association signals to each SNP, SNP-window p-values were treated as a SNP prioritization

192    metric by generating the geometric mean of $-\log_{10}$(p-values) across windows containing each

193    SNP.  That is, the prioritization metric was calculated using the p-value for each SKAT analysis

194    window ($p_i$) that contained the $k^{th}$ SNP ($n$ analysis windows). With the exception of the first and

195    last five SNPs in a region, each $SNP_k$ was included in five analysis windows (n=5). Thus, for

196    each SNP $k$, we calculated its prioritization metric as:

$$\text{Prioritization Metric } SNP_k = -\log_{10}\left(\prod_{i=1}^{n} p_i\right)^{\frac{1}{n}}$$

197                                                                                          (Equation 1)

198          *Bayesian Approach: Credible SNP Sets*.  Frequentist approaches, such as those

199    implemented SKAT or single-SNP logistic regression analyses are widely utilized; however,

200    their resulting p-values are not without limitations.(51) For one, p-values do not capture the

201    confidence of a particular association. Furthermore, they're more dependent on factors such as

202    the power of the statistical test (influenced by sample size and other variables). Bayesian

203    methods offer an alternative approach; here, Bayes factors are used, capturing the ratio of

204    probabilities between the null and alternative hypotheses.

205          As a comparison to the frequentist approaches, we used SNPTEST to generate the

206    Bayes factors (BF), using the score test and additive genotype modeling.(52) Posterior

207    probabilities for a given SNP $k$, were then calculated using method published by the Welcome

9

208    Trust Case Control Consortium.(53) For SNPs 1-j in the region, the posterior probability for each

209    SNP *k,* was calculated by:

$$\text{Posterior Probability for SNP}_k = \frac{\text{BF}_k}{\Sigma_j \, \text{BF}_j}$$

210    (Equation 2)

211    Using these posterior probabilities, the 95% credible set was determined for each region. This

212    test assumes only one causal SNP in the region and places equal *a priori* probabilities that the

213    causal SNP is any one of the analyzed SNPs.(53)  In this study, we applied this method to

214    previously defined regions (43) where we hypothesized the association signal is driven by one

215    SNP.

216         Like the single-SNP logistic regression analyses, this Bayesian analysis is not weighted

217    by functional data. Thus, for a ΔMGW-weighted analysis, a derived credible set was generated

218    from posterior probabilities that accounted for each SNP's ΔMGW through *ad hoc* weighting,

219    where the posterior probability for a given SNP *k,* was calculated by weighting the Bayes factor

220    by $\Delta\text{MGW}_k$ divided by the weighted average of Bayes factors for SNPs 1-j in the region. Here,

221    the derived posterior probability for each SNP *k*, is:

$$\text{Derived Posterior Probability for SNP}_k = \frac{\text{BF}_k \, \Delta\text{MGW}_k}{\Sigma_j \, \text{BF}_j \, \Delta\text{MGW}_j}$$

222    (Equation 3)

223    Using these values, the derived 95% credible SNP sets were generated and compared with the

224    unweighted 95% credible SNP sets. This methodology enabled weighting by a continuous

225    variable versus existing methods designed for dichotomous (presence/absence of functional

226    annotation) SNP weights.(54)

227    **Functional Annotation**

10

228        To evaluate the functional plausibility for an identified variant, several publically available

229   resources were referenced. For variant associations with gene expression (eQTL status), the

230   Genotype-Tissue Expression (GTEx) dataset, version 7 (hg19) was queried at

231   gtexportal.org.(55) GTEx is a comprehensive eQTL resource, providing eQTL information

232   across 48 tissues. SNPs were also queried using the SCREEN (Search Candidate cis-

233   Regulatory Elements by Encode, http://screen.encodeproject.org).(56, 57) Built using Encode

234   data, SCREEN (hg19) evaluates if a given genomic coordinate resides in a Candidate cis-

235   Regulatory Element (ccRE). ccREs are designated based on evidence from DNase

236   hypersensitivity sites, H3K4me3 and H3K27ac histone activity, and CTCF-binding data.

237   SCREEN contains 1.31 million ccREs, correlating to 20.8% of the mappable human genome

238   (http://screen.encodeproject.org). Genomic variants were also evaluated for evidence of long-

239   range DNA interaction via Hi-C data (hg19) available through the Yue Lab 3D Genome Browser

240   (http://promoter.bx.psu.edu/hi-c/).(58) Similar to the ccRE search, SNPs were queried to see if

241   they resided in a genome region that exhibited long-range chromatin interactions. The Yue

242   Lab's Capture Hi-C data offers information across 19 cell line options. We evaluated immune-

243   related cell types: naïve B-Cells, CD4_Total (CD4 activated and Naïve), CD8 naïve, monocytes,

244   and neutrophils.

# 245  Results
246

247  ***For ∆MGW, SNPs in the human genome exhibit a stronger right skewed distribution in***

248  ***comparison to the complete sample space.***

249        In the complete sample space of ∆MGW, ∆MGW values ranged from 0.00 to 3.16 Å, with

250   a mean of 0.77 Å and a standard deviation of 0.42. **(Table 1)** The overall data exhibited a right-

251   skewed distribution (**Figure 3**) with few sequences inducing large changes in MGW.

252   Unsurprisingly, given the sequence-dependency of this topological measure, parsing the data

253    by the paired alleles (fifth nucleotide, see Methods and Materials), revealed allele-specific

254    patterns of ∆MGW **(Table 1)**.  Transition pairings (A/G and C/T) yielded the smallest changes in

255    ∆MGW, while transversion pairings (Purine/Pyrimidine) produced the largest changes in ∆MGW.

256    Subsets that represent complimentary allele pairs (i.e. A/G & T/C; A/C & T/G) yielded the same

257    ∆MGW values. (**Table 1**) Of all allele-pairings, A/T alleles presented the largest ∆MGW with a

258    mean of 1.16 Å  (SD, 0.47) **(Figure 3)**.

259         We compared the ∆MGW sample space statistics to the observed frequencies of ∆MGW

260    across the human genome using dbSNP data. The hg19 download of NCBI dbSNP150

261    contained 234,104,110 entries. After pruning to high quality (see Methods and Materials), bi-

262    allelic SNPs, 199,038,197 polymorphisms remained. For these SNPs, there was an average

263    ∆MGW of 0.68 Å with a standard deviation of 0.43. In comparison to the ∆MGW sample space,

264    SNPs across the genome exhibited a stronger, right-skewed distribution of ∆MGW. (**Figure 3**).

265    Transition SNPs are more likely to occur (59, 60), and this is consistent with our SNP150

266    summarizations, where transition SNPs comprised 66.43% of the dataset (**Table S1**). Our

267    ∆MGW sample space summarization showed that transition allele pairings had the smallest

268    change in ∆MGW **(Table 1);** thus, the decreased average in ∆MGW dbSNP data is expected

269    and illustrates the high prevalence shape-preserving SNPs in the genome. To evaluate patterns

270    in ∆MGW by SNP function (i.e. missense, intron, coding-synonymous), SNPs with a single

271    NCBI-designation (see Methods and Materials) were subset and summarized (**Table 2, Figure**

272    **4**). Notably, some SNP categories are limited to specific sequence combinations(61) (i.e. stop-

273    loss, **Table S2**)**,** which were reflected in the SNP-function-specific patterns of ∆MGW. (**Figure**

274    **4**) Coding-synonymous SNPs exhibited the smallest overall change in ∆MGW (mean=0.48 Å).

275    Unknown and intron SNPs, which are not constrained to specific sequences (by definition),

276    comprised the two largest categories ($n_{unknown}$=99,004,130; $n_{intron}$=84,909,115) and yielded high

277    averages for ∆MGW: 0.69 Å and 0.56 Å, respectively.

278   ***Fine-mapping SLE-associated genomic regions using ΔMGW prioritization identifies***

279   ***potentially functional SNPs.***

280   To-date, more than 100 genomic loci have been associated with SLE.(43, 62) Here, we

281   selected the genomic regions containing *FAM167A-BLK*, *STAT4*, and *TNIP1* for fine-mapping

282   because these regions showed robust single-SNP associations ($p < 5 \times 10^{-8}$) with SLE in at least

283   two ancestries (*FAM167A-BLK*: EA and AA; *STAT4*: EA and HA; *TNIP1*: EA and HA) and the

284   association signals are not refined to a single SNP, due in part to strong linkage disequilibrium.

285   Furthermore, neither the SNPs nor their LD proxies are protein-coding variants, leaving DNA

286   topology as a potential functional mechanism. For each region, we first describe the previous

287   SNP association results (43) and their LD patterns, by ancestry. Each region is then

288   summarized by its ΔMGW measures which were used in frequentist and Bayesian ΔMGW-

289   weighted analyses. SNPs identified by the ΔMGW-weighted analyses were subsequently

290   investigated for existing functional evidence (See Methods and Materials).

291   ***FAM167A-BLK.***

292   The SLE-associated region at 8p23 lies upstream of *FAM167A* and *BLK*, which are in a

293   head-to-head gene orientation. Across the 500kb candidate region, 835 and 933 genotyped

294   SNPs passed quality control in the EA and AA data, respectively. In the previous(43) logistic

295   regression analyses, the primary peak of association was captured by a 60 kb window. In EA,

296   the most significant SNP associations mapped to a 26 kb region of 16 SNPs in high LD ($r^2 > 0.8$);

297   within the AA data, the top associations were refined to a smaller 14 kb window containing 7

298   highly correlated SNPs (**Figure 5**).The summary statistics for ΔMGW for SNPs in the 500 kb

299   and 60 kb regions were comparable to what was observed across the genome, with only a few

300   SNPs imposing large changes in MGW (**Table S3**).

13

301     Hypothesizing that plausibly functional SNPs can be identified by incorporating both

302     ΔMGW and evidence for disease association, we applied two ΔMGW-weighted approaches via

303     SKAT and Bayesian credible sets. For the 500 kb region, SKAT was applied in a 5-SNP rolling

304     window (see Methods and Materials). Across the region, SNPs with the highest SKAT-weighted

305     prioritizations largely followed the pattern observed in the single-SNP logistic regression

306     analyses. That is, SNPs that were not previously associated with SLE were not prioritized solely

307     on ΔMGW, as illustrated in the region outside of the 40 kb peak of association (**Figure 5**). When

308     weighted by ΔMGW, rs2061831 was sharply prioritized in both the EA and AA analyses (**Figure**

309     **5**). In EA, rs2061831 was one of the 14 highly correlated SNPs identified by the single-SNP

310     logistic regression analyses; likewise, in AA, it was also within the LD block comprising the 7

311     most highly associated SNPs. While the other SNPs in these LD blocks exhibited comparable

312     SLE-association, rs2061831 had the greatest ΔMGW at 1.63 Å, prioritizing it above other SNPs

313     in the weighted analyses. Importantly, while the single-SNP logistic regression analyses

314     identified a different top SNP in EA (rs13277113) and AA (rs2736440), ΔMGW-weighting

315     prioritized the same SNP (rs2061831), across ancestries. An unweighted SKAT prioritized the

316     signal downstream of rs2061831, to the region where multiple SNPs from the same highly-

317     associated LD block were included in the same 5-SNP windows (**Figure S1, Tables S4-S5**).

318     The ΔMGW-weighted frequentist fine-mapping evidence for rs2061831 was

319     corroborated using the Bayesian refinement approach.  In both EA and AA, the derived ΔMGW-

320     weighted credible set placed the highest posterior probability on rs2061831 (58.9%-EA; 44.2%-

321     AA) (**Figure 5**). In the un-weighted (standard) Bayesian analysis, rs2061831 was included in the

322     EA (30.6% posterior probability) and AA (20.9% posterior probability) 95% credible sets, but it

323     was not the highest prioritized **(Table S4-S5)**. Instead, the SNPs originally identified in the

324     ancestry-specific logistic regression analyses were given the highest posterior probability—EA:

325     rs13277113 (49.9% posterior probability), AA: rs2736340 (33.1%).  Thus, like the frequentist

326    approach, weighting by ΔMGW resolved the signal in both EA and AA to the same SNP,

327    rs2061831.

328          Using ΔMGW as a prioritization metric, rs2061831 was consistently prioritized in both

329    EA and AA data. SNP rs2061831 has a ΔMGW of 1.63 Å, which is 2 standard deviations above

330    the mean across dbSNP150. Interestingly, this SNP is a transition polymorphism

331    (Purine/Purine), a polymorphism type which we previously showed to have the smallest (on

332    average) ΔMGW (**Table 1, Figure 3**). Considering only transition SNPs, rs2061831 is actually

333    4.52 standard deviations above the mean $\Delta MGW_{\text{transition SNPs}}$ (0.50 Å), indicating a considerable

334    departure from the expected value and thus we would hypothesize a greater likelihood of

335    functional relevance. Given the consistent evidence for a signal at rs2061831 in both the EA

336    and AA data, we explored previously described (see Methods and Materials) functional data

337    resources for evidence of biological relevance, in comparison to the top SNP signals from the

338    single-SNP analyses (rs13277113 in EA; and rs2736440 in AA). All three SNPs are in high LD

339    ($R^2$>0.95) with one another in both EUR and AFR 1000 genomes data. Thus, it is unsurprising

340    that all three SNPs yielded similar eQTL results via GTEx (data not shown). Despite the high

341    LD, these three SNPs are physically separated by several kilobases. Of these three SNPs,

342    rs2061831 is the only SNP that maps (via SCREEN) to a Candidate Cis-Regulatory Element

343    (accession number: EH37E0941109) showing evidence for DNase, H3K27ac, and CTCF-

344    binding activity. Consulting the 3D-genome browser yielded a larger number of long-range

345    chromatin interactions in monocytes, B-Cells, and CD4 cells for rs2061831, in comparison to

346    rs13277113 and rs2736440 (**Figure S2**). Thus, in this region, ΔMGW-weighting successfully

347    differentiated among highly-correlated SNPs and prioritized rs2061831, a SNP within a

348    potentially important regulatory region as documented by independent data.

349    ***STAT4***

15

350    The single-SNP SLE associations at 2q32 span the *STAT4* gene **(Figure 6)**. SNP

351    associations reached genome significance in the EA and HA cohorts, with the strongest signals

352    within intronic regions**.**(43) In the 500 kb region, there were 192 and 202 genotyped SNPs that

353    passed quality control measures in EA and HA, respectively. In both ancestries, the primary

354    peak of association was captured by a broad 110 kb window (**Figure 6**). The strongest

355    associations in the EA data (p-values < $1\times10^{-62}$) mapped to six SNPs in high LD, spanning 29

356    kb. Five of these SNPs also comprised the LD block of strongest associations in the HA data

357    (p< $1\times10^{-13}$), in a slightly narrower 26 kb region. The consistency of SNP association results in

358    the EA and HA data provided a prime opportunity to test ΔMGW-prioritization among highly-

359    correlated SNPs.

360    The mean ΔMGW for SNPs in this region was 0.72 Å in EA and 0.73 Å in HA and both

361    cohorts had a median ΔMGW of 0.56 Å. While these average ΔMGW were slightly higher than

362    what was observed across the entire bi-allelic dbSNP dataset (mean=0.68 Å), the EA and HA

363    medians were of the same magnitude (dbSNP ΔMGW median=0.56). The ΔMGW for SNPs

364    within the 110 kb association window exhibited similar means as the 500 kb region (**Table S6**).

365    We again applied the two ΔMGW-weighted approaches using SKAT and Bayesian

366    credible sets in the region. In EA, the ΔMGW-weighted SKAT analyses shifted the top signal

367    upstream to rs11889341, which markedly increased its priority (**Figure 6**). This SNP was one of

368    the top six SNPs in the single-SNP association LD-block. While it and the other five SNPs were

369    all significantly associated with SLE, rs11889341 had the greatest ΔMGW at 1.75 Å, which

370    prioritized it over the other SNPs in the LD block; the remaining SNPs had ΔMGW values

371    ranging from 0.31-1.12 Å (**Figure 6**). In HA, weighting by ΔMGW in the SKAT analysis also

372    prioritized rs11889341 as the top SNP. This SNP was previously identified with the best p-value

373    in the single-SNP association analysis, but in the ΔMGW-weighted approach, its prioritization

374    distinctly increased relative to the other SNPs in the LD block (**Figure 6**).

16

375       In the Bayesian analysis, rs11889341 was included in the EA and HA derived ΔMGW-

376    weighted 95% credible sets (**Figure 6**). In EA, rs11889341 was not in the unweighted 95%

377    credible set but inclusion of ΔMGW increased its posterior probability from 2.4% to 6.0% (**Table**

378    **S7, Figure S3**). In EA, rs7568275 yielded the strongest signal in both the unweighted (81.0%

379    posterior probability) and derived ΔMGW-weighted (77.3% posterior probability) credible sets

380    (**Table S7**). This is important to note, as rs7568275 had a much smaller ΔMGW (0.66 Å) than

381    rs11889341 (1.75 Å.). This provided an example where the magnitude of the Bayes factor was

382    so large ($p=4 \times 10^{68}$), that the influence of ΔMGW was largely diminished in the analysis.

383    However, despite the predominant rs7568275 signal, the derived credible set still detected

384    rs11889341, the SNP identified by the ΔMGW-weighted SKAT approach. In the HA data,

385    rs11889341 yielded the largest posterior probability in the ΔMGW-weighted derived credible set.

386    This SNP also had the largest posterior probability in the unweighted credible set. Unlike the EA

387    analysis, where the magnitude of the Bayes factor dominated the impact of the ΔMGW-

388    weighting, in the HA data, the ΔMGW strongly increased the posterior probability of rs11889341

389    from 58.6% to 73.5% (**Figure 6, Table S8**). This limited the derived 95% credible set to only 3

390    SNPs: rs11889341 (73.5%), rs8179673 (16.6%), and rs7574865 (4.8%) (**Table S8**).

391       In the single-SNP association analyses of *STAT4* SNPs, the association signal was

392    refined to an LD block of 6 SNPs in the EA data and 5 SNPs in the HA dataset. In ΔMGW-

393    weighted analyses, rs11889341 was sharply prioritized over other SNPs in the LD block, with an

394    exception in the EA ΔMGW-weighted derived credible set, where the high magnitude of the

395    Bayes factor for rs7568275 ($bf=2.20 \times 10^{64}$) over other SNPs ($bf \leq 1.79 \times 10^{63}$) largely negated

396    any impact of ΔMGW in this analysis. Considering the evidence for rs11889341 in the other

397    three analyses due to its strong combination of SLE association and ΔMGW, we would

398    hypothesize that rs11889341 would be a candidate functional polymorphism. Like rs2061831 in

399    *FAM167A-BLK*, rs11889341 is also a transition SNP (purine/purine). While transition SNPs are

400    more frequent across the genome (previously shown in Table S1), there are few transition SNPs

401    (+/- 4 nucleotides) that yield such a high ΔMGW (mean ΔMGW for transition SNPs=0.50 Å).

402    Evaluation of publically available functional datasets (see METHODS) yielded limited

403    information for both rs7568275 and rs11889341. Neither of these SNPs were identified as

404    eQTLs in GTEx nor were they within Candidate Cis-Regulatory regions (cCREs). Furthermore,

405    neither variant was shown with long range chromatin interactions in the in the currently available

406    HI-C data via the 3D genome browser. However, despite the lack of functional information from

407    these resources, functional evaluation of rs11889341 is available via a 2018 study by Patel and

408    colleagues, where transancestral mapping identified rs11889341 with strong association with

409    SLE.[63] In this study, rs11889341 was associated with *STAT1* expression in B-cells through

410    increased binding of the transcription factor, HMGA1. Given the relationship between

411    transcription factor binding and DNA topology[20, 31, 32, 64, 65], we hypothesize that the

412    identified functional activity of rs11889341 (via HMGA1 binding) may be mediated by the large

413    MGW change imposed by the SNP's alleles.

414    **TNIP1**

415        Previous single-SNP association analyses[43] identified genome-wide significant

416    findings ($p<5x10^{-8}$) in EA and HA data at 5q33 (**Figure 7**).  In the 500 kb region, there were 497

417    and 500 high quality genotyped SNPs in the EA and HA data, respectively. The peak of SLE

418    association is captured by a 40 kb window which encompasses most of the *TNIP1* gene. In the

419    EA data, the top associations mapped to three SNPs (rs960709, rs10036748, rs6889239) in

420    high LD, spanning 3 kb of a *TNIP1* intron. These three SNPs are also encompassed by the

421    associated LD block in the HA data, where four, highly correlated SNPs (rs1422673, rs960709,

422    rs10036748, and rs6889239) yielded p-values $< 5x10^{-8}$. As completed in the *FAM167A-BLK* and

423    *STAT4* regions, we again applied ΔMGW-weighted fine-mapping strategies to prioritize these

424    non-coding SLE-associated SNPs.

425    In the *TNIP3* region, the lists of high-quality genotyped SNPs were largely the same

426    between the EA and HA datasets. Consequently, the statistics for ΔMGW in this region were

427    very similar between the two cohorts. Across the 500 kb window of high quality SNPs, the

428    average ΔMGW was 0.67 Å (median=0.55 Å) in both EA and HA. (**Table S9**) These values

429    were slightly lower than the observed mean for bi-allelic SNPs from dbSNP (**Table 1**).

430    The SKAT analyses yielded similar results between the EA and HA data. The ΔMGW-

431    weighted analyses did not effectively prioritize or refine the SNP signal. Unlike *FAM167A-BLK*

432    and *STAT4*, ΔMGW-weighting did not resolve the top signal to the same SNP in both

433    ancestries. Instead, in *TNIP1*, the top SNPs in the ΔMGW-weighted analyses for EA

434    (rs6889239) and HA (rs10036748) were the same as those identified in the single-SNP logistic

435    regression analysis (**Figure 7**).  The SNPs that were prioritized in the unweighted SKAT

436    analyses were also prioritized in the ΔMGW-weighted analyses; notably, in this region ΔMGW-

437    weighting actually dampened the signal because the SNPs with the greatest SLE association

438    values had low magnitudes of ΔMGW (ranging from 0.31-0.37 Å). This pattern was also

439    observed in the Bayesian approach, where SNPs with the highest posterior probabilities in the

440    derived credible sets exhibited lower posterior probabilities than in the unweighted credible set

441    **(Figures 7 and S4 and Tables S10-S11**), again due to the low magnitudes of ΔMGW for top-

442    associated SNPs.

443    In *TNIP1*, the ΔMGW-weighted analyses did not differentially prioritize SNPs in

444    comparison to the unweighted approaches. While there were SNPs with large ΔMGW in the

445    region, these did not have strong SLE-associations. Unlike the *FAM167A-BLK* and *STAT4*

446    regions, where ΔMGW successfully prioritized specific SNPs, this was not achieved in the

447    *TNIP1* region. This could indicate several possibilities, including:  ΔMGW may not be a relevant

448    mechanism for these SNPs, another DNA measure may be more informative, DNA topology

449    may not be a functional driver for this region, and/or or the functional variant was not included in

450    these analyses. Here, an alternative strategy is required to identify the most plausible functional

451    polymorphisms.

452    **Discussion**

453          Sequence-dependent DNA topology could provide important functional context for

454    associations, especially for polymorphisms that do not impose protein changes (e.g., coding-

455    synonymous) and/or variants mapping to non-coding regions.  We explored ΔMGW, a specific

456    sequence-dependent measure of DNA topology, as a weighting variable in fine-mapping

457    analyses. In a sample of 300k SNPs, Wang *et al*. previously found that MGW-preserving SNPs

458    are more common.(42) Here, we built upon these findings through a full census of bi-allelic

459    SNPs (n=199,038,197) across the genome. We showed the observed genomic ΔMGW was

460    significantly lower than the complete ΔMGW sample space. These findings were consistent with

461    the relative frequencies of transversion (~33%) and transition (~66%) mutations in the human

462    genome.(59, 60) We hypothesized that phenotypically-associated SNPs with large ΔMGW

463    would be more likely to impose functional consequences; and thus, proposed ΔMGW as a

464    prioritization metric in fine-mapping studies.

465          We tested our hypothesis using ΔMGW weights in two fine-mapping approaches in three

466    regions (*FAM167A-BLK*, *STAT4*, and *TNIP1*) with well-established SLE associations. In

467    *FAM167A-BLK and STAT4*, we successfully identified SNPs of possible functional

468    consequence, underscoring ΔMGW as a plausibly informative prioritization metric in fine-

469    mapping studies.

470          There are several advantages to using sequence dependent topology, such as ΔMGW, as a

471    weighting metric in fine-mapping studies. For one, it is an intrinsic variable, inherent to the

472    genetic sequence surrounding the polymorphism; thus, it is not reliant on external data which

473    may offer limited information for the SNPs of interest (database bias). As an intrinsic variable it

474    is also not ancestry specific, tissue specific, or sample size dependent.  Limitations in external

475   (non-intrinsic) data may down-weight potentially causal SNPs due to a lack of available

476   functional data. While publically available functional resources continue to expand, they still

477   present these challenges, especially for rare or novel variants. This is particularly relevant for

478   diverse study populations where annotation resources based on European data offer

479   inadequate or no coverage for regions of interest.(14) For example, Sherman *et al.* presented

480   deep sequencing in 910 individuals of African descent and found over 296 million base pairs

481   which were absent in the human reference genome.(15) Novel variants or regions are unlikely

482   to be annotated by commonly used resources. Therefore, while a SNP's functional relevance

483   can be supported by public resources, a lack of information does not necessarily indicate a

484   variant's lack of function. This was illustrated by rs11889341 in *STAT4*, which lacked functional

485   information from public resources (GTEx, ENCODE, 3D-genome browser)(55, 56, 58), but in a

486   targeted functional study by Patel *et al.*, rs11889341 was correlated with gene expression and

487   binding of the transcription factor HMGA1.(63)  We identified rs11889341 using ΔMGW as the

488   prioritizing variable. Thus, prioritizing SNPs by a factor intrinsic to DNA may help alleviate some

489   bias that would otherwise be introduced by missing data from publically available functional

490   datasets. Consequently, we propose including ΔMGW among annotation resources used in

491   SNP-weighted fine-mapping methods.

492      Changes in DNA topology can potentially impact an array of biological functions such as

493   transcription factor binding, chromatin remodeling, or methylation.(20, 21, 23, 26, 31, 32, 36)

494   Likewise, using DNA topology as a SNP prioritization metric does not limit functional information

495   to a single biological mechanism. This may be especially beneficial when the relationship

496   between phenotype and biological mechanism is unknown. While functional work in *STAT4*

497   showed that rs11889341 altered HMGA1 binding, functional work is still needed to evaluate the

498   rs2061831 genotype in *FAM167A-BLK*. Here, the biological implications of rs2061831 could

499   involve transcription factor binding, and/or, given its apparent location within a long-range

500    chromatin interaction hotspot (**Figure S1**), chromatin organization. Considering the strong trans-

501    ancestral signal of rs2061831 across EA and AA, further functional work should explore whether

502    this SNP acts through an independent functional mechanism or through interactions with other

503    variants in the region (e.g. within the context of sequence-dependent structural motifs), such as

504    the insertion-deletion identified in a study of ATAC-seq data in 100 individuals of British

505    Ancestry.(66)  Leveraging changes in DNA topology can identify potentially causal

506    polymorphisms and also generate specific hypotheses for functional follow-up studies.

507    Furthermore, sequence-dependent DNA topology is a weighting scheme that informatively

508    decouples SNPs in high LD, a long sought after feature as associations and eQTLs are often

509    confounded by LD. In *FAM167A-BLK*, we observed comparable eQTL evidence for SNPs in the

510    associated LD cluster, making eQTL status ineffective at differentiating highly-correlated SNPs.

511    Instead, consideration of sequence-dependent ΔMGW allowed differential prioritization among

512    these otherwise, highly-correlated SNPS, selecting rs2061831 as a plausible functional

513    candidate SNP.

514        Another advantage to using local DNA topology in fine-mapping studies is its consistency of

515    information across ancestries. Assuming identical flanking sequence (e.g., no genomic variant

516    within +/- 4 bases of the SNP), a SNP's impact on DNA topology would be constant across

517    ancestries, highlighting the potential utility of DNA topology as a means of resolving association

518    signals across ancestries. Here, we showed that ΔMGW-weighted analyses of *FAM167A-BLK*

519    and *STAT4* resolved the association signal to the same SNP in each ancestry via the frequentist

520    approach, followed by largely corroborating evidence via the derived credible sets in the

521    Bayesian approach. Notably, rs2061831 was not the top-associated SNP in either the ancestry-

522    specific analyses; however, it was previously identified via the SLE Immunochip trans-ancestral

523    meta-analysis, where combining association signals across ancestries identified it as the top

524    SNP.(43)

525 **Limitations and Future Work**

526      There are several considerations and limitations to using sequence-dependent topology

527 as a weighting metric in fine-mapping analyses. Notably, some of these limitations could result

528 in inconclusive and/or insignificant results, as observed in the *TNIP1* region.  First, the functional

529 variants may not have been genotyped or imputed in the study. Analyses that utilize SNP-

530 specific weights decouple associations from LD. Thus, a weighted metric performs best when

531 the functional SNP is included in the analysis set.  For this reason, we propose application of

532 this prioritization technique in genomic regions where there is high confidence that the functional

533 variants have been genotyped or imputed. We note this limitation exists for any statistical

534 association method.

535      Second, DNA topology, here ΔMGW, may not be the mechanism impacting phenotype.

536 While sequence dependent DNA topology can influence a number of functional factors(18, 21,

537 23, 24, 32),  it is not the only source of biological interactions and could be irrelevant for a

538 specific phenotype. Thus, when using change in DNA topology, such as ΔMGW, in fine-

539 mapping studies, analyses should be considered in the form of a two-parameter hypothesis – a

540 combination of association signal and ΔMGW. For example, in both the *FAM167A-BLK* and

541 *STAT4* regions, the highest prioritized SNPs, rs2061831 and rs11889341, did not have the

542 largest magnitude of ΔMGW in the regions (**Figures 5-6**). Instead, these two SNPs were

543 prioritized by their combined SLE-association and ΔMGW.

544      Third, we placed greater weights on SNPs with larger magnitudes of change on DNA

545 topology. We recognize that even small changes could yield functional consequences. Thus,

546 future studies should explore weighting SNPs by particular topological profiles (e.g., those

547 matching binding site profiles). For instance, our *TNIP1* analyses did not show strong signals

548 when weighting by the magnitude of ΔMGW, but this does not definitively rule out MGW as a

549 functional mechanism (e.g. driven by pattern, not magnitude). The focus on MGW was

550    motivated by the breadth of study on MGW and function.(18, 20, 32, 34, 36) So while this

551    manuscript considered a single parameter, ΔMGW, we are currently expanding to incorporate

552    additional measures (e.g., helix twist, roll) through multivariate approaches that account for the

553    correlation structure (dependencies) among spatial measures.

554       Fourth, in this study, we used SKAT and a derived credible sets (Bayesian) approach to

555    apply a topological weighting scheme to prioritize SNPs; however, we note that there are other

556    methods that can incorporate weights for SNP association analyses.(10, 67)  Here, we assumed

557    that the majority of variants in the region are non-causal, which is why we selected SKAT over a

558    combined burden test. However, we note that the results from SKAT and SKAT-O were largely

559    similar. Similarly, in case of the Bayesian approach applied here, a limitation is its assumption

560    that a single causal SNP exists in a region, but other Bayesian methods can be explored.(53,

561    68)  In the EA *STAT4* data, the magnitudes of the Bayes factors were so large that weighting by

562    ΔMGW yielded minimal impact. Future work should consider approaches to scale weighting

563    schemes by a constant derived from the magnitude of signal across a genomic region. In the

564    SKAT approach, for the sliding analysis window, we used five SNPs, which should yield a

565    region that is neither too wide nor too unstable. Additional testing could potentially improve

566    optimization of parameters for this analysis. Furthermore, we emphasize that our evaluation of

567    the SKAT results by summarizing each SNP as the geometric mean of SKAT-analysis p-values

568    should be regarded as a metric for prioritizing SNPs, not an association analyses, as these

569    values do not have the statistical properties of a p-value. Overall, these limitations should be

570    carefully considered when applying these specific methods; but they also highlight opportunities

571    to further explore the relationship between sequence-dependent DNA topology and phenotype

572    associations.

573       In summary, weighting SNP associations by functional data can greatly improve

574    identification of potentially causal SNPs; however, existing annotation resources can negatively

575 affect these outcomes when SNP information is unavailable in public datasets, especially in

576 non-EA populations.(8, 10, 11, 14)  In this study, we presented and tested sequence-dependent

577 DNA topology as a novel annotation source for genetic fine-mapping studies. As an intrinsic

578 property, sequence-dependent DNA shape alleviates many of the challenges imposed by

579 external data resources; and it provides potential functional (testable) context for associations

580 (e.g. topological disruption for protein binding). Using ΔMGW in weighted analyses, we

581 successfully prioritized functional SNPs in two SLE-associated regions with high LD. Likewise,

582 as an annotation resource, sequence-dependent DNA topology, such as ΔMGW, is readily

583 applicable in any fine-mapping methods that can incorporate continuous values for SNP

584 weights. Altogether, this manuscript presents methods that are immediately applicable to

585 existing genetic data, and it illustrates how sequence-dependent DNA topology can be used as

586 a paradigm to investigate and understand genetic associations in fine-mapping studies.

587

591

592 **Declaration of Interests.**

593 The authors declare no competing interests.

594

## References.

599

600 1. MacArthur,J., Bowler,E., Cerezo,M., Gil,L., Hall,P., Hastings,E., Junkins,H., McMahon,A.,
601     Milano,A., Morales,J., *et al.* (2017) The new NHGRI-EBI Catalog of published genome-
602     wide association studies (GWAS Catalog). *Nucleic Acids Res.*, **45**, D896–D901.

603 2. Visscher,P.M., Brown,M.A., McCarthy,M.I. and Yang,J. (2012) Five years of GWAS
604     discovery. *Am. J. Hum. Genet.*, **90**, 7–24.

605 3. Visscher,P.M., Wray,N.R., Zhang,Q., Sklar,P., McCarthy,M.I., Brown,M.A. and Yang,J. (2017)
606     10 Years of GWAS Discovery: Biology, Function, and Translation. *Am. J. Hum. Genet.*,
607     **101**, 5–22.

608 4. McCarthy,M.I., Abecasis,G.R., Cardon,L.R., Goldstein,D.B., Little,J., Ioannidis,J.P.A. and
609     Hirschhorn,J.N. (2008) Genome-wide association studies for complex traits: consensus,
610     uncertainty and challenges. *Nat. Rev. Genet.*, **9**, 356–369.

611 5. Manolio,T.A., Collins,F.S., Cox,N.J., Goldstein,D.B., Hindorff,L.A., Hunter,D.J.,
612     McCarthy,M.I., Ramos,E.M., Cardon,L.R., Chakravarti,A., *et al.* (2009) Finding the
613     missing heritability of complex diseases. *Nature*, **461**, 747–753.

614 6. Pasaniuc,B. and Price,A.L. (2017) Dissecting the genetics of complex traits using summary
615     association statistics. *Nat. Rev. Genet.*, **18**, 117–127.

616 7. Farh,K.K.-H., Marson,A., Zhu,J., Kleinewietfeld,M., Housley,W.J., Beik,S., Shoresh,N.,
617     Whitton,H., Ryan,R.J.H., Shishkin,A.A., *et al.* (2015) Genetic and epigenetic fine
618     mapping of causal autoimmune disease variants. *Nature*, **518**, 337–343.

619 8. Gomez-Cabrero,D., Abugessaisa,I., Maier,D., Teschendorff,A., Merkenschlager,M., Gisel,A.,
620     Ballestar,E., Bongcam-Rudloff,E., Conesa,A. and Tegnér,J. (2014) Data integration in
621     the era of omics: current and future challenges. *BMC Syst. Biol.*, **8**, I1.

622 9. Faye,L.L., Machiela,M.J., Kraft,P., Bull,S.B. and Sun,L. (2013) Re-Ranking Sequencing
623     Variants in the Post-GWAS Era for Accurate Causal Variant Identification. *PLOS Genet.*,
624     **9**, e1003609.

625 10. Kichaev,G., Yang,W.-Y., Lindstrom,S., Hormozdiari,F., Eskin,E., Price,A.L., Kraft,P. and
626     Pasaniuc,B. (2014) Integrating Functional Data to Prioritize Causal Variants in Statistical
627     Fine-Mapping Studies. *PLOS Genet.*, **10**, e1004722.

628 11. Xu,Z. and Taylor,J.A. (2009) SNPinfo: integrating GWAS and candidate gene information
629     into functional SNP selection for genetic association studies. *Nucleic Acids Res.*, **37**,
630     W600–W605.

631 12. Lee,S., Wu,M.C. and Lin,X. (2012) Optimal tests for rare variant effects in sequencing
632     association studies. *Biostatistics*, **13**, 762–775.

633 13. Nicolae,D.L., Gamazon,E., Zhang,W., Duan,S., Dolan,M.E. and Cox,N.J. (2010) Trait-
634     Associated SNPs Are More Likely to Be eQTLs: Annotation to Enhance Discovery from
635     GWAS. *PLOS Genet.*, **6**, e1000888.

636   14. Kessler,M.D., Yerges-Armstrong,L., Taub,M.A., Shetty,A.C., Maloney,K., Jeng,L.J.B.,
637       Ruczinski,I., Levin,A.M., Williams,L.K., Beaty,T.H., *et al.* (2016) Challenges and
638       disparities in the application of personalized genomic medicine to populations with
639       African ancestry. *Nat. Commun.*, **7**, 12521.

640   15. Sherman,R.M., Forman,J., Antonescu,V., Puiu,D., Daya,M., Rafaels,N., Boorgula,M.P.,
641       Chavan,S., Vergara,C., Ortega,V.E., *et al.* (2019) Assembly of a pan-genome from deep
642       sequencing of 910 humans of African descent. *Nat. Genet.*, **51**, 30–35.

643   16. Need,A.C. and Goldstein,D.B. (2009) Next generation disparities in human genomics:
644       concerns and remedies. *Trends Genet. TIG*, **25**, 489–494.

645   17. Manrai,A.K., Funke,B.H., Rehm,H.L., Olesen,M.S., Maron,B.A., Szolovits,P.,
646       Margulies,D.M., Loscalzo,J. and Kohane,I.S. (2016) Genetic Misdiagnoses and the
647       Potential for Health Disparities. *N. Engl. J. Med.*, **375**, 655–665.

648   18. Privalov,P.L., Dragan,A.I., Crane-Robinson,C., Breslauer,K.J., Remeta,D.P. and
649       Minetti,C.A.S.A. (2007) What Drives Proteins into the Major or Minor Grooves of DNA?
650       *J. Mol. Biol.*, **365**, 1–9.

651   19. Yakovchuk,P., Protozanova,E. and Frank-Kamenetskii,M.D. (2006) Base-stacking and
652       base-pairing contributions into thermal stability of the DNA double helix. *Nucleic Acids*
653       *Res.*, **34**, 564–574.

654   20. Yang,L., Orenstein,Y., Jolma,A., Yin,Y., Taipale,J., Shamir,R. and Rohs,R. (2017)
655       Transcription factor family-specific DNA shape readout revealed by quantitative
656       specificity models. *Mol. Syst. Biol.*, **13**, 910.

657   21. Duan,C., Huan,Q., Chen,X., Wu,S., Carey,L.B., He,X. and Qian,W. (2018) Reduced intrinsic
658       DNA curvature leads to increased mutation rate. *Genome Biol.*, **19**, 132.

659   22. Sati,S. and Cavalli,G. (2017) Chromosome conformation capture technologies and their
660       impact in understanding genome function. *Chromosoma*, **126**, 33–44.

661   23. Lazarovici,A., Zhou,T., Shafer,A., Dantas Machado,A.C., Riley,T.R., Sandstrom,R.,
662       Sabo,P.J., Lu,Y., Rohs,R., Stamatoyannopoulos,J.A., *et al.* (2013) Probing DNA shape
663       and methylation state on a genomic scale with DNase I. *Proc. Natl. Acad. Sci. U. S. A.*,
664       **110**, 6376–6381.

665   24. Abe,N., Dror,I., Yang,L., Slattery,M., Zhou,T., Bussemaker,H.J., Rohs,R. and Mann,R.S.
666       (2015) Deconvolving the recognition of DNA shape from sequence. *Cell*, **161**, 307–318.

667   25. Bansal,M., Kumar,A. and Yella,V.R. (2014) Role of DNA sequence based structural features
668       of promoters in transcription initiation and gene expression. *Curr. Opin. Struct. Biol.*, **25**,
669       77–85.

670   26. Parker,S. and Tullius,T.D. (2011) DNA shape, genetic codes, and evolution. *Curr. Opin.*
671       *Struct. Biol.*

672   27. Olson,W.K., Bansal,M., Burley,S.K., Dickerson,R.E., Gerstein,M., Harvey,S.C.,
673       Heinemann,U., Lu,X.-J., Neidle,S., Shakked,Z., *et al.* (2001) A Standard Reference

Frame for the Description of Nucleic Acid Base-pair Geometry. *J. Mol. Biol.*, **313**, 229–237.

28. Lu,X.-J. and Olson,W.K. (1999) Resolving the discrepancies among nucleic acid conformational analyses11Edited by I. Tinoco. *J. Mol. Biol.*, **285**, 1563–1575.

29. Dickerson,R.E. (1989) Definitions and nomenclature of nucleic acid structure components. *Nucleic Acids Res.*, **17**, 1797–1803.

30. Rohs,R., West,S.M., Sosinsky,A., Liu,P., Mann,R.S. and Honig,B. (2009) The role of DNA shape in protein-DNA recognition. *Nature*, **461**, 1248–1253.

31. Meysman,P., Marchal,K. and Engelen,K. (2012) DNA structural properties in the classification of genomic transcription regulation elements. *Bioinforma. Biol. Insights*, **6**, 155–168.

32. Stella,S., Cascio,D. and Johnson,R.C. (2010) The shape of the DNA minor groove directs binding by the DNA-bending protein Fis. *Genes Dev.*, **24**, 814–826.

33. Irobalieva,R.N., Fogg,J.M., Catanese,D.J., Catanese,D.J., Sutthibutpong,T., Chen,M., Barker,A.K., Ludtke,S.J., Harris,S.A., Schmid,M.F., *et al.* (2015) Structural diversity of supercoiled DNA. *Nat. Commun.*, **6**, 8440.

34. Morgunova,E., Yin,Y., Jolma,A., Dave,K., Schmierer,B., Popov,A., Eremina,N., Nilsson,L. and Taipale,J. (2015) Structural insights into the DNA-binding specificity of E2F family transcription factors. *Nat. Commun.*, **6**, 10050.

35. Ngo,T.T.M., Zhang,Q., Zhou,R., Yodh,J.G. and Ha,T. (2015) Asymmetric unwrapping of nucleosomes under tension directed by DNA local flexibility. *Cell*, **160**, 1135–1144.

36. Perino,M., van Mierlo,G., Karemaker,I.D., van Genesen,S., Vermeulen,M., Marks,H., van Heeringen,S.J. and Veenstra,G.J.C. (2018) MTF2 recruits Polycomb Repressive Complex 2 by helical-shape-selective DNA binding. *Nat. Genet.*, **50**, 1002–1010.

37. Chen,C. and Pettitt,B.M. (2016) DNA Shape versus Sequence Variations in the Protein Binding Process. *Biophys. J.*, **110**, 534–544.

38. Shepherd,J.W., Greenall,R.J., Probert,M.I.J., Noy,A. and Leake,M.C. (2020) The emergence of sequence-dependent structural motifs in stretched, torsionally constrained DNA. *Nucleic Acids Res.*, 10.1093/nar/gkz1227.

39. Chiu,T.-P., Comoglio,F., Zhou,T., Yang,L., Paro,R. and Rohs,R. (2016) DNAshapeR: an R/Bioconductor package for DNA shape prediction and feature encoding. *Bioinformatics*, **32**, 1211–1213.

40. Zhou,T., Shen,N., Yang,L., Abe,N., Horton,J., Mann,R.S., Bussemaker,H.J., Gordân,R. and Rohs,R. (2015) Quantitative modeling of transcription factor binding specificities using DNA shape. *Proc. Natl. Acad. Sci.*, **112**, 4654–4659.

41. Duzdevich,D., Redding,S. and Greene,E. (2014) DNA Dynamics and Single-Molecule Biology. *Chem. Rev.*, **114**, 3072–3086.

711  42. Wang,X., Zhou,T., Wunderlich,Z., Maurano,M.T., DePace,A.H., Nuzhdin,S.V. and Rohs,R.
712      (2018) Analysis of Genetic Variation Indicates DNA Shape Involvement in Purifying
713      Selection. *Mol. Biol. Evol.*, **35**, 1958–1967.

714  43. Langefeld,C.D., Ainsworth,H.C., Graham,D.S.C., Kelly,J.A., Comeau,M.E., Marion,M.C.,
715      Howard,T.D., Ramos,P.S., Croker,J.A., Morris,D.L., *et al.* (2017) Transancestral
716      mapping and genetic load in systemic lupus erythematosus. *Nat. Commun.*, **8**, 16021.

717  44. van Dijk,M. and Bonvin,A.M.J.J. (2009) 3D-DART: a DNA structure modelling server.
718      *Nucleic Acids Res.*, **37**, W235-239.

719  45. Pettersen,E.F., Goddard,T.D., Huang,C.C., Couch,G.S., Greenblatt,D.M., Meng,E.C. and
720      Ferrin,T.E. (2004) UCSF Chimera--a visualization system for exploratory research and
721      analysis. *J. Comput. Chem.*, **25**, 1605–1612.

722  46. Haeussler,M., Zweig,A.S., Tyner,C., Speir,M.L., Rosenbloom,K.R., Raney,B.J., Lee,C.M.,
723      Lee,B.T., Hinrichs,A.S., Gonzalez,J.N., *et al.* (2019) The UCSC Genome Browser
724      database: 2019 update. *Nucleic Acids Res.*, **47**, D853–D858.

725  47. Karolchik,D., Hinrichs,A.S., Furey,T.S., Roskin,K.M., Sugnet,C.W., Haussler,D. and
726      Kent,W.J. (2004) The UCSC Table Browser data retrieval tool. *Nucleic Acids Res.*, **32**,
727      D493–D496.

728  48. Lander,E.S., Linton,L.M., Birren,B., Nusbaum,C., Zody,M.C., Baldwin,J., Devon,K.,
729      Dewar,K., Doyle,M., FitzHugh,W., *et al.* (2001) Initial sequencing and analysis of the
730      human genome. *Nature*, **409**, 860–921.

731  49. Wu,M.C., Lee,S., Cai,T., Li,Y., Boehnke,M. and Lin,X. (2011) Rare-variant association
732      testing for sequencing data with the sequence kernel association test. *Am. J. Hum.
733      Genet.*, **89**, 82–93.

734  50. Lee,S., Emond,M.J., Bamshad,M.J., Barnes,K.C., Rieder,M.J., Nickerson,D.A., NHLBI GO
735      Exome Sequencing Project—ESP Lung Project Team, Christiani,D.C., Wurfel,M.M. and
736      Lin,X. (2012) Optimal unified approach for rare-variant association testing with
737      application to small-sample case-control whole-exome sequencing studies. *Am. J. Hum.
738      Genet.*, **91**, 224–237.

739  51. Stephens,M. and Balding,D.J. (2009) Bayesian statistical methods for genetic association
740      studies. *Nat. Rev. Genet.*, **10**, 681–690.

741  52. Marchini,J., Howie,B., Myers,S., McVean,G. and Donnelly,P. (2007) A new multipoint
742      method for genome-wide association studies by imputation of genotypes. *Nat. Genet.*,
743      **39**, 906–913.

744  53. The Wellcome Trust Case Control Consortium, Maller,J.B., McVean,G., Byrnes,J.,
745      Vukcevic,D., Palin,K., Su,Z., Howson,J.M.M., Auton,A., Myers,S., *et al.* (2012) Bayesian
746      refinement of association signals for 14 loci in 3 common diseases. *Nat. Genet.*, **44**,
747      1294–1301.

748  54. Kichaev,G., Roytman,M., Johnson,R., Eskin,E., Lindström,S., Kraft,P. and Pasaniuc,B.
749      (2017) Improved methods for multi-trait fine mapping of pleiotropic risk loci. *Bioinforma.*
750      *Oxf. Engl.*, **33**, 248–255.

751  55. GTEx Consortium (2015) The Genotype-Tissue Expression (GTEx) pilot analysis:
752      Multitissue gene regulation in humans. *Science*, **348**, 648–660.

753  56. ENCODE Consortium (2004) The ENCODE (ENCyclopedia Of DNA Elements) Project.
754      *Science*, **306**, 636–640.

755  57. ENCODE Project Consortium (2012) An integrated encyclopedia of DNA elements in the
756      human genome. *Nature*, **489**, 57–74.

757  58. Wang,Y., Song,F., Zhang,B., Zhang,L., Xu,J., Kuang,D., Li,D., Choudhary,M.N.K., Li,Y.,
758      Hu,M., *et al.* (2018) The 3D Genome Browser: a web-based browser for visualizing 3D
759      genome organization and long-range chromatin interactions. *Genome Biol.*, **19**, 151.

760  59. Nachman,M.W. and Crowell,S.L. (2000) Estimate of the mutation rate per nucleotide in
761      humans. *Genetics*, **156**, 297–304.

762  60. Zhao,Z. and Boerwinkle,E. (2002) Neighboring-Nucleotide Effects on Single Nucleotide
763      Polymorphisms: A Study of 2.6 Million Polymorphisms Across the Human Genome.
764      *Genome Res.*, **12**, 1679–1686.

765  61. Kitts,A., Phan,L., Ward,M. and Holmes,J.B. (2014) The Database of Short Genetic Variation
766      (dbSNP) National Center for Biotechnology Information (US).

767  62. Niewold,T.B. (2015) Advances in Lupus Genetics. *Curr. Opin. Rheumatol.*, **27**, 440–447.

768  63. Patel,Z.H., Lu,X., Miller,D., Forney,C.R., Lee,J., Lynch,A., Schroeder,C., Parks,L.,
769      Magnusen,A.F., Chen,X., *et al.* (2018) A plausibly causal functional lupus-associated
770      risk variant in the STAT1-STAT4 locus. *Hum. Mol. Genet.*, **27**, 2392–2404.

771  64. Parvin,J.D. and Sharp,P.A. (1993) DNA topology and a minimal set of basal factors for
772      transcription by RNA polymerase II. *Cell*, **73**, 533–540.

773  65. Scaffidi,P. and Bianchi,M.E. (2001) Spatially Precise DNA Bending Is an Essential Activity of
774      the Sox2 Transcription Factor. *J. Biol. Chem.*, **276**, 47296–47302.

775  66. Kumasaka,N., Knights,A.J. and Gaffney,D.J. (2019) High-resolution genetic mapping of
776      putative causal interactions between regions of open chromatin. *Nat. Genet.*, **51**, 128–
777      137.

778  67. Yang,J., Fritsche,L.G., Zhou,X. and Abecasis,G. (2017) A Scalable Bayesian Method for
779      Integrating Functional Information in Genome-wide Association Studies. *Am. J. Hum.*
780      *Genet.*, **101**, 404–416.

781  68. Jiang,J., Cole,J.B., Freebern,E., Da,Y., VanRaden,P.M. and Ma,L. (2019) Functional
782      annotation and Bayesian fine-mapping reveals candidate genes for important agronomic
783      traits in Holstein bulls. *Commun. Biol.*, **2**, 212.
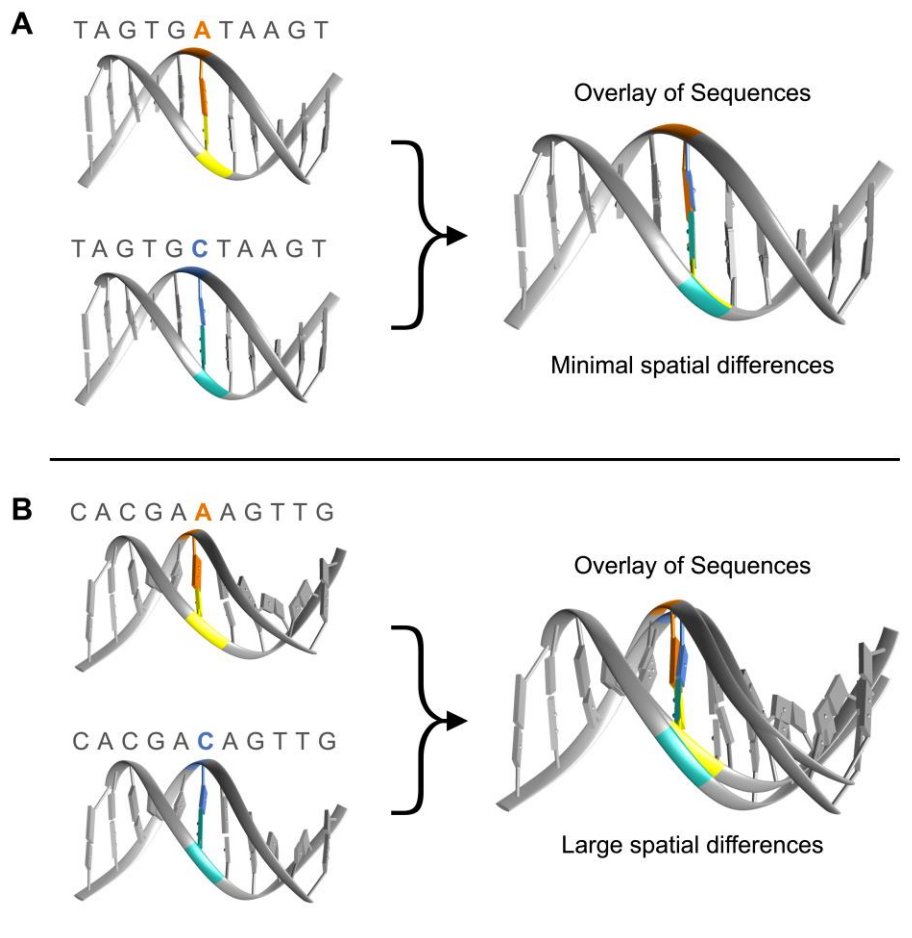
784

785    **Primary Figures and Legends**

786

787    **Figure 1: Single nucleotide substitutions sequence can impose large or small changes**

788    **on local DNA shape, dependent on the flanking sequence.**

789    (A) A single A/C substitution within a sequence generates minimal spatial differences.

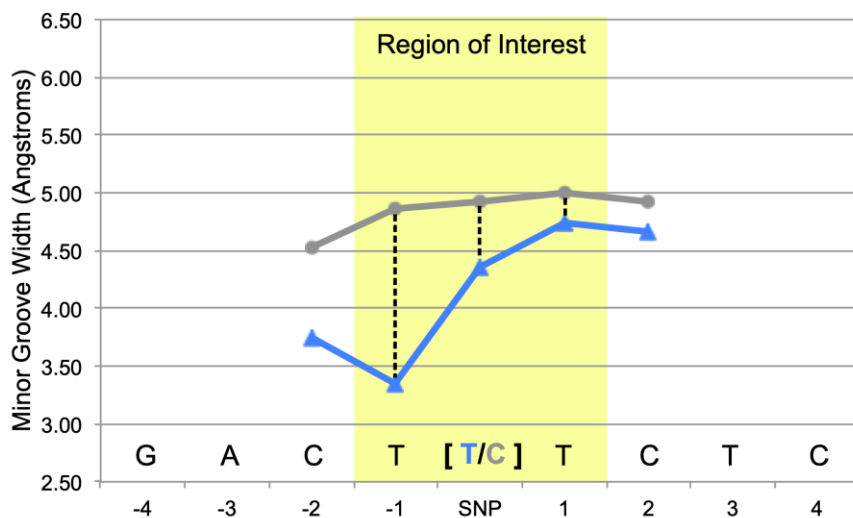790    (B) A single A/C substitution within a sequence imposes large spatial differences



791

792  **Figure 2: Generation of ∆MGW for a SNP.**

793  (A) Minor groove width measures are plotted for the two sequences generated by a specific bi-

794  allelic T/C SNP. For a given SNP, the flanking sequence (+/- 4 bp) was used as input for

795  DNAshapeR (via Bioconductor) which calculates MGW along a rolling sequence window.  For a

796  9-mer sequence, the MGW can be consistently provided at the SNP's position +/- one

797  nucleotide which is highlighted in yellow and labeled as the 'region of interest'. Expanding this

798  region to additional nucleotides would require a longer input sequence and increases chance of

799  additional variants being within the input (and introducing additional variability). Although the two

800  sequences for a SNP only differ at one nucleotide (at the SNP position), the impact on MGW

801  carries through adjacent bases. Thus, ∆MGW was calculated to capture the change in MGW for

802  a SNP by incorporating information at the SNP's position and +/- 1 base pair (dashed lines).

803  (B) Workflow for calculating the ∆MGW for a bi-allelic SNP. This method captures the change in

804  MGW at the SNP position and +/- 1 base pair. This Euclidean distance captures ∆MGW as a

805  measure of magnitude (in Angstroms).

**A** MGW for a bi-allelic (T/C) SNP



**B** Calculation of ΔMGW for a bi-allelic SNP

T/C SNP with four flanking nucleotides (+/-)
G-A-C-T-[T/C]-T-C-T-C

Two unique sequences per bi-allelic SNP
G-A-C-T-**C**-T-C-T-C
G-A-C-T-**T**-T-C-T-C

Generate MGW prediction using DNAshapeR for SNP and +/- 1 bp (k=-1,0,1)

$k=-1$      $k=0$      $k=1$

3.35  –  **4.36**  –  4.74
4.86  –  **4.93**  –  5.00

Calculate Euclidean distance (ΔMGW)

$$\sqrt{\sum_{k=-1}^{1}(MGW_k^{allele1} - MGW_k^{allele2})^2} = 1.63$$

806

807

808 **Figure 3: Summarization of ΔMGW across the complete sample space**

809 (A) ΔMGW sample space was constructed on six allele pairings (A/C, A/G, A/T, C/G, C/T, G/T)

810 with all possible combinations for flanking +/- 4 bp. This yielded 393,216 paired sequences that

811 were evaluated for ΔMGW.

812 (B) The distribution of ΔMGW for the 393,216 paired sequences, these summary statistics are

813 listed in Table 1.

814 (C) Two randomly selected paired sequences from the average and right tail of the ΔMGW

815 distribution are shown. Sequences are plotted with their respective MGW values (Angstroms).

816 ΔMGW is calculated as a Euclidean distance, which captures the change in MGW (dashed

817 lines) at the SNP position and +/- 1bp (highlighted in orange). ATGA[C/A]CGAT exhibits a small

818 ΔMGW , at 0.47 Å while TCCA[T/A]ATTG yields a large change in MGW (2.34 Å) which we

819 would hypothesize to have greater potential for functional consequence if also associated with
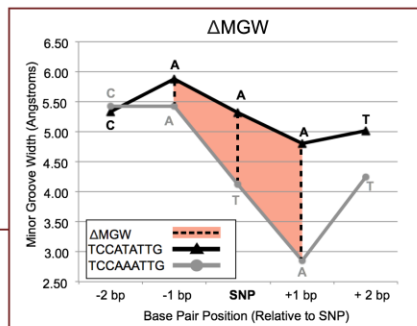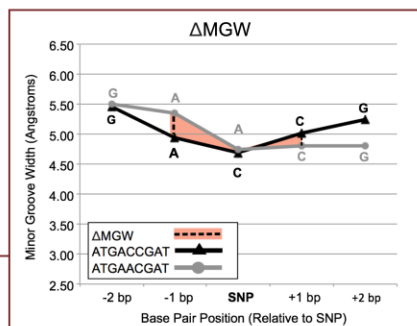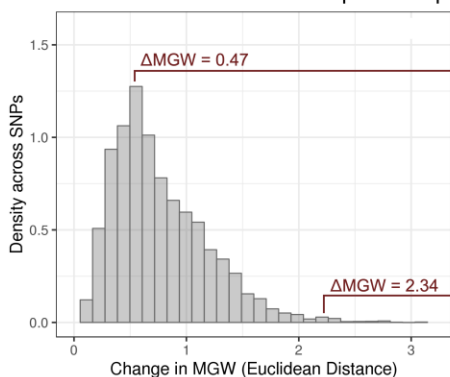
820 disease status.

821 (D) The ΔMGW distribution for all paired sequences (gray) is shown superimposed on the

822 ΔMGW distributions by 5th nucleotide alleles (blue). Transition pairings (C/T, A/G) have a more

823 strongly skewed distribution with a smaller average ΔMGW compared to transversion pairings

824 (A/C, A/T, C/G, G/T), (Table 1). Pairings that represent complimentary sequences (C/T – A/G

825 and A/C – T/G) exhibit the same distributions of ΔMGW, as expected.

826

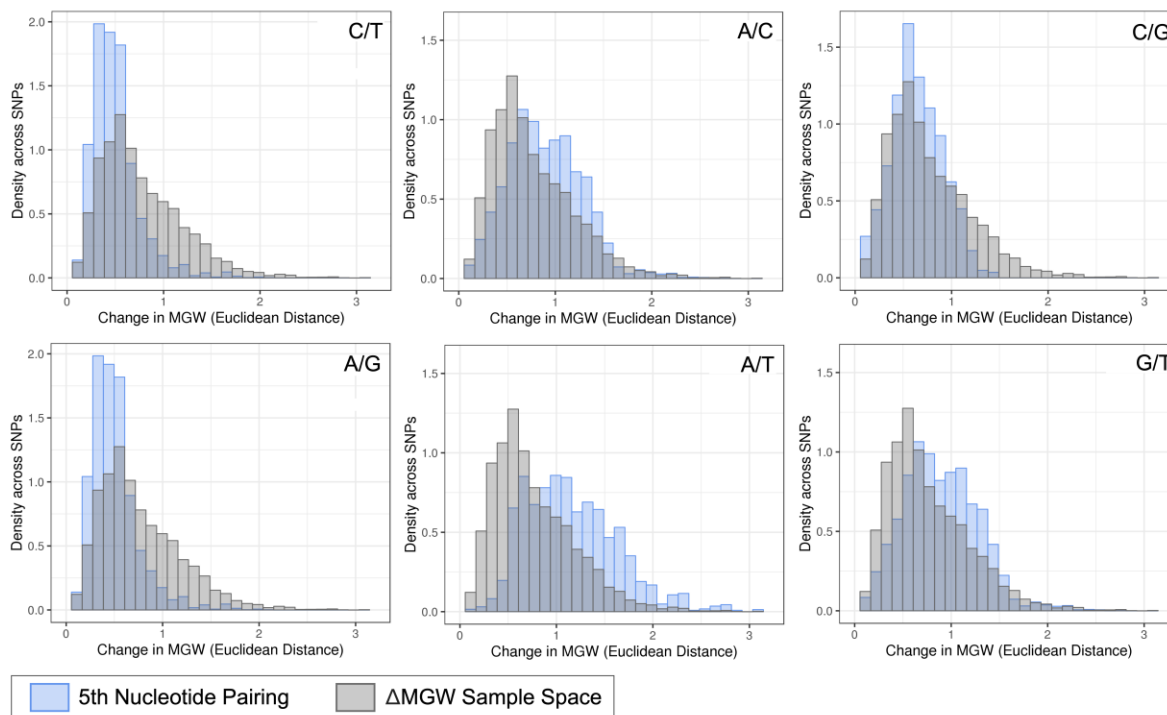**A** Generation of all possible sequence pairings (with +/- 4 bp)

**C**

| Sequence combinations by paired alleles | Total |
|---|---|
| (4)-(4)-(4)-(4)-**[A/C]**-(4)-(4)-(4)-(4) | = 65,536 |
| (4)-(4)-(4)-(4)-**[A/G]**-(4)-(4)-(4)-(4) | = 65,536 |
| (4)-(4)-(4)-(4)-**[A/T]**-(4)-(4)-(4)-(4) | = 65,536 |
| (4)-(4)-(4)-(4)-**[C/G]**-(4)-(4)-(4)-(4) | = 65,536 |
| (4)-(4)-(4)-(4)-**[C/T]**-(4)-(4)-(4)-(4) | = 65,536 |
| (4)-(4)-(4)-(4)-**[G/T]**-(4)-(4)-(4)-(4) | = 65,536 |
| | Total=393,216 |

**B** Distribution of ΔMGW across complete sample space



**D**

Transition Pairings
(Purine/Purine or Pyrimidine/Pryimidine)

Transversion Pairings
(Purine/Pyrimidine)


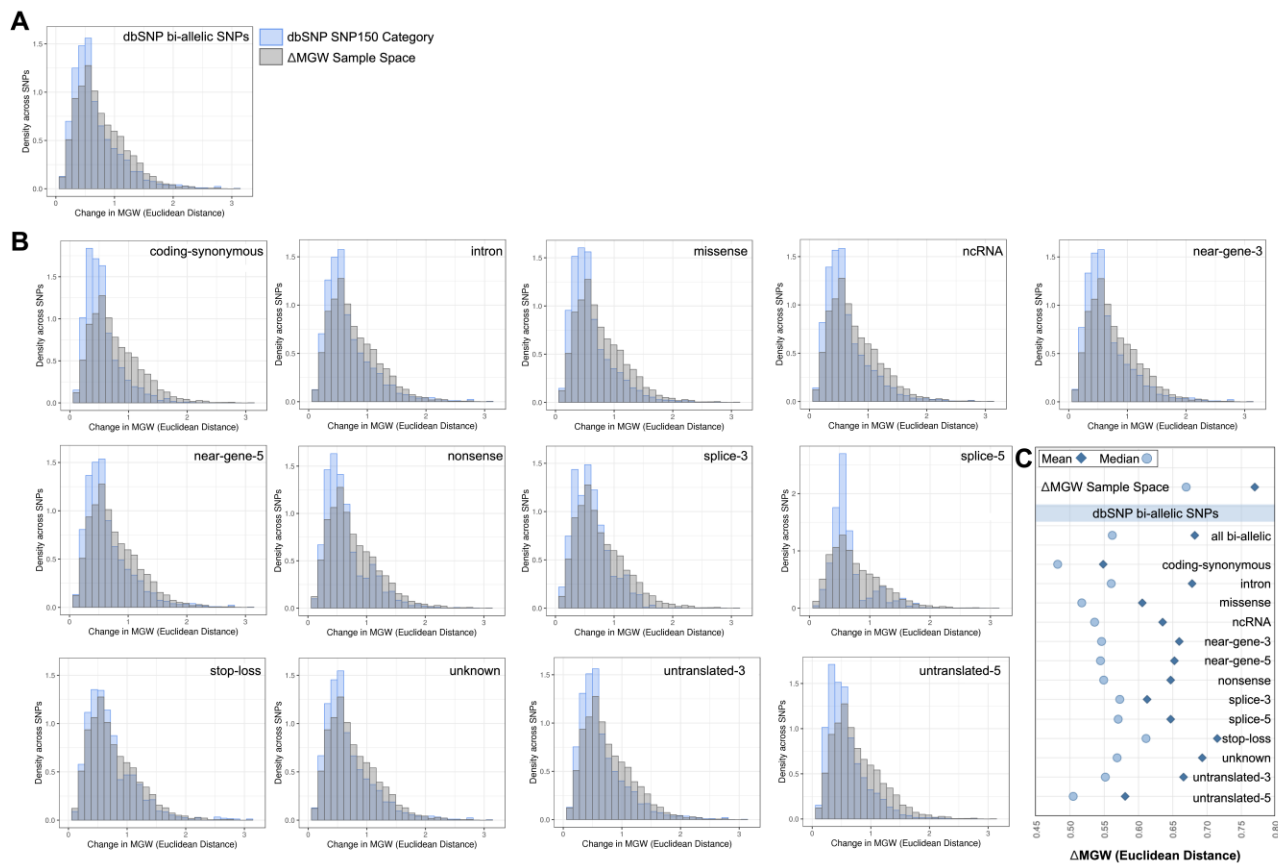
5th Nucleotide Pairing   ΔMGW Sample Space

827

828

35

829 **Figure 4: Summarization of ∆MGW across the human genome using bi-allelic SNPs from**

830 **dbSNP SNP150.**

831 (A) Comparison of ∆MGW sample space (Figure 3) and the observed ∆MGW from SNPs across

832 the genome (via dbSNP). Distribution of ∆MGW is shown in blue for observed bi-allelic SNPs

833 from the SNP150 dataset (n=199,038,197 SNPs). The ∆MGW sample space distribution (Figure

834 3) is plotted in gray (n=393,216 paired sequences). The observed ∆MGW across genomic

835 SNPs showed a stronger right skewed distribution than what would be expected from a random

836 sampling of the entire sample space of all-possible sequences.  Only small numbers of SNPs

837 elicit large magnitudes of ∆MGW.

838 (B) ∆MGW distributions are similarly shown for SNP subsets, by NCBI function (exclusive NCBI

839 function label for each SNP, see Methods and Materials). Again, each distribution is

840 superimposed with the distribution from the ∆MGW sample space (shown in gray). Subsetting

841 by NCBI function yields similar patterns observed in part A, with observed genomic SNPs

842 showing smaller averages in ∆MGW. Some NCBI SNP-functions have specific sequence

843 requirements (Supplemental Table 1) and these are reflected in the resulting ∆MGW

844 distributions which are also sequence-dependent (e.g. splice-6, nonsense).

845 (C) The mean and median ∆MGW for each SNP category. All dbSNP SNP categories have

846 significantly lower mean and median compared to the ∆MGW sample space (Tables 1-2).

847 Coding-synonymous SNPs have the smallest magnitudes of ∆MGW, compared to all other
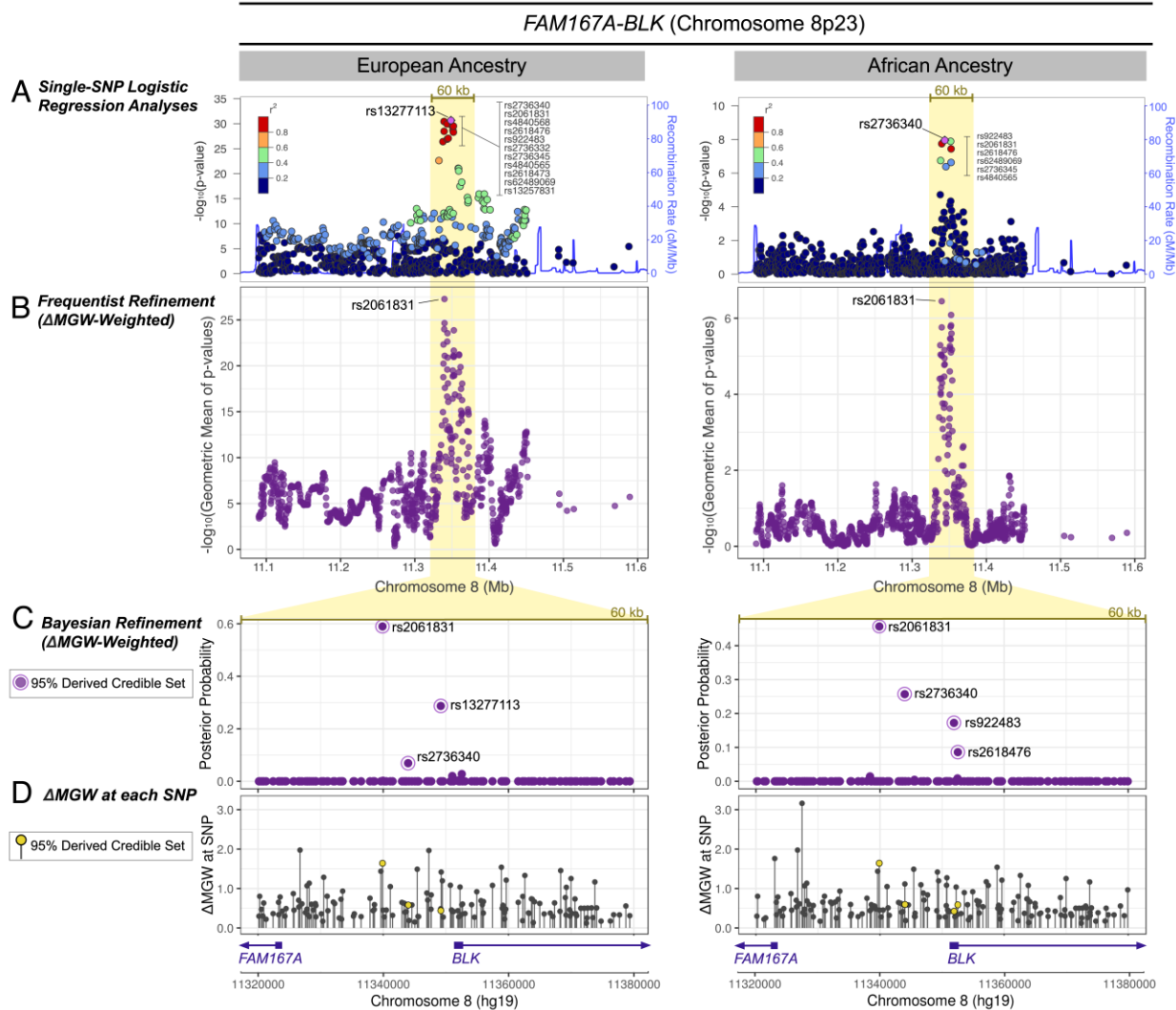
848 categories.

849

850

851

852

853 **Figure 5: *FAM167A-BLK* ΔMGW prioritization by Frequentist and Bayesian Methods in**
854 **European and African Ancestries.**

855 (A) Genotyped SNPs that passed quality control and were within 250kb of the top single-SNP

856 association analysis in EA and AA data. A 60 kb region capturing the primary peak of

857 association is highlighted. In both the EA and AA data a cluster of SNPs in high LD yielded the

858 top association signals.

859 (B) Using SKAT as a ΔMGW-weighted frequentist approach, rs2061831 was sharply prioritized

860 over SNPs in the previously identified LD blocks. While the single-SNP logistic regression

861 analyses in (A) identified a different top SNP in the EA (rs13277113) and AA (rs2736340) data,

862 rs2061831 was consistently prioritized as the top SNP in both the EA and AA analyses. ΔMGW-

863 weighting did not yield spurious associations for with SNPs outside the broad 60 kb peak of

864 association highlighted in yellow.

865 (C) SNPs within the 60 kb association peak were analyze by a Bayesian approach. The ΔMGW-

866 weighted posterior probabilities are plotted. While the majority of SNPs yielded infinitesimal

867 posterior probabilities, those comprising the 95% derived credible sets are labeled. Akin to the

868 ΔMGW-weighted SKAT analyses, rs2061831 was again prioritized in both the EA and the AA

869 data, with the largest posterior probability.

870 (D) The ΔMGW is plotted for each SNP in the 60 kb region. The ΔMGW for a SNP is sequence-

871 specific thus yielding the same values in EA and AA data. Differences between the two plots

872 result from differences in genotyped SNP lists (i.e. SNPs that are monomorphic in one

873 population would not be plotted). SNPs identified by the derived ΔMGW-weighted credible set

874 are plotted in yellow. While rs2061831 had a large ΔMGW, other SNPs in the region had larger

875 magnitudes of ΔMGW but did not show evidence of SLE-association. This illustrates the 2-

876 parameter hypothesis of considering a combination of association signal and magnitude of

877 ΔMGW. Prioritized SNPs fall upstream of both *FAM167A* and *BLK*.

38

878

879    **Figure 6: STAT4 ΔMGW prioritization by Frequentist and Bayesian Methods in European**
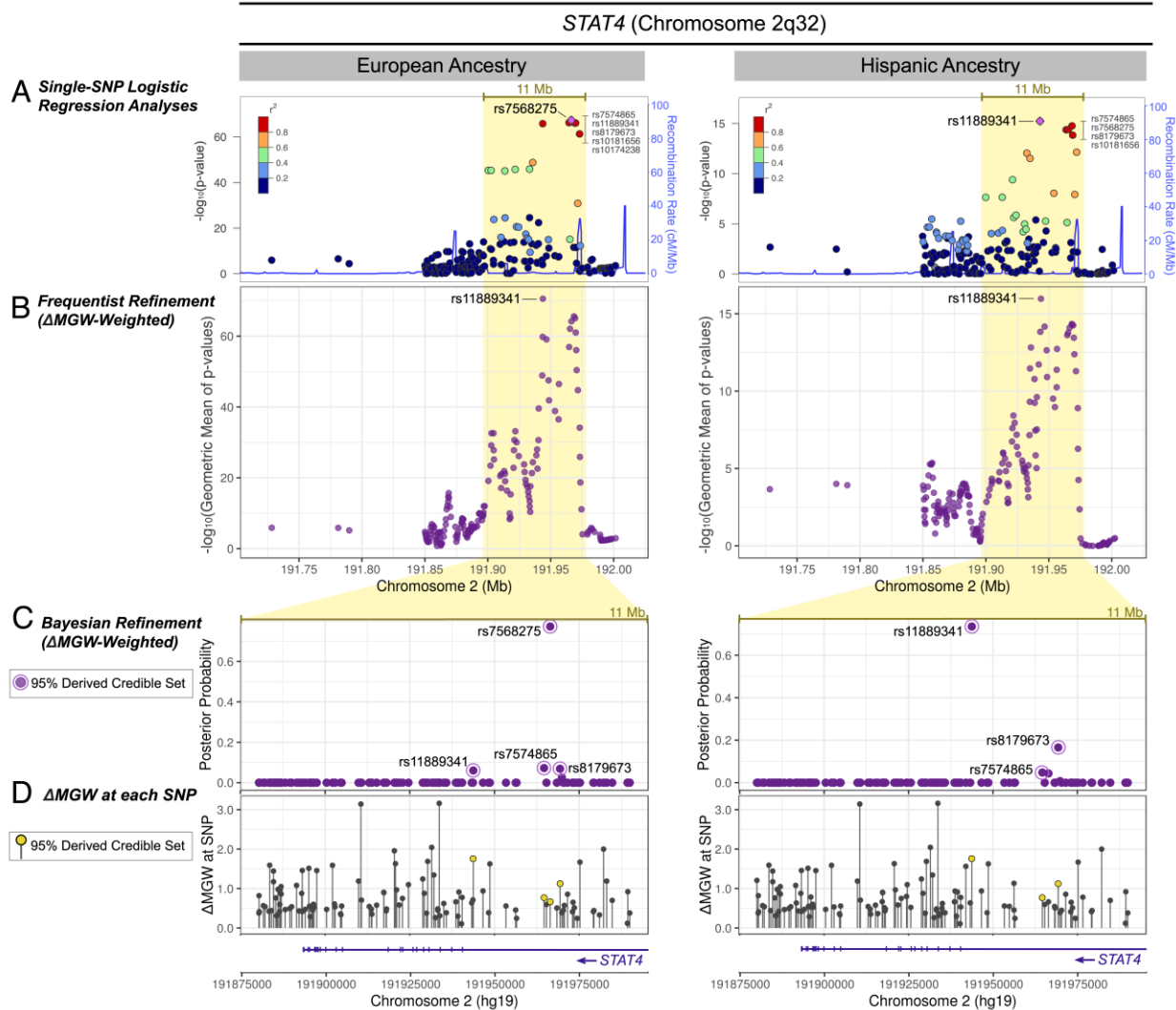
880    **and Hispanic Ancestries.**

881    (A) Regional association plots in EA and HA for genotyped SNPs that passed quality control

882    and were within 250kb of the top single-SNP association analysis in *STAT4.* Within the broad 11

883    Mb peak of association (highlighted in yellow), a cluster of SNPs in high LD yielded the top

884    association values.

885    (B) SNP refinement using SKAT with a ΔMGW-weighted approach sharply prioritizes

886    rs11889341 in both EA and HA data. In the EA data, the ΔMGW-weighting shifted the top signal

887    to rs1188931, whereas in the HA data, it simply further accentuated the signal above other

888    SNPs.

889    (C) For the highlighted 11 Mb region, SNP posterior probabilities are plotted for the derived,

890    ΔMGW-weighted Bayesian analysis. While the frequentist MGW-weighted approach prioritized

891    the same SNP (rs1188931) in both ancestries, this was not observed in the Bayesian approach.

892    In the EA data, the Bayes factor for rs7568275 (BF=2.20x10$^{64}$) was at such a large magnitude,

893    that it was largely unaffected by ΔMGW-weighting. However, rs1188931 still entered the 95%

894    derived credible set, but with a much smaller posterior probability (6.03%) compared to

895    rs7568275 (77.25%). In the HA data, ΔMGW-weighting further prioritized rs1188931.

896    (D) The ΔMGW for SNPs within the 11 Mb region. SNPs that were identified by the derived

897    ΔMGW-weighted credible set are plotted in yellow. Again, the analytic approaches consider

898    SNPs in the context of a 2-parameter hypothesis, evaluating SNPs for a combination of

899    association signal and magnitude of ΔMGW. Hence, the prioritized SNPs (yellow) are not

900    necessarily the SNPs with the largest ΔMGW in the region. Prioritized SNPs occur within an

901    intron of *STAT4.*

902

903

41

904    **Figure 7: *TNIP1* ∆MGW prioritization by Frequentist and Bayesian Methods in European**

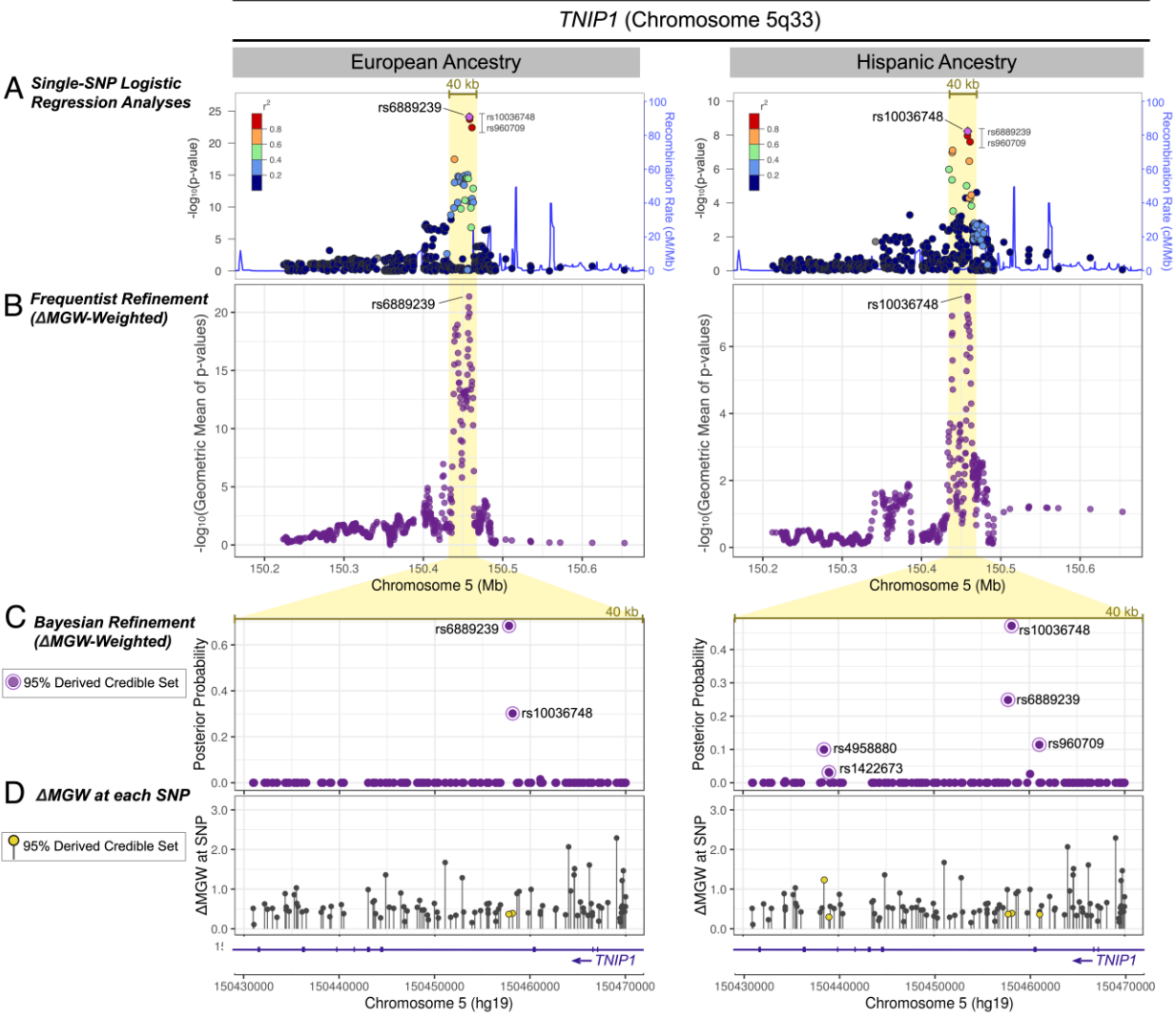905    **and Hispanic Ancestries.**

906    (A) Genotyped SNPs within 250 kb of the top single-SNP association analysis are shown for EA

907    and HA. The 40 kb region that captures the primary peak of association is highlighted in yellow.

908    In EA and HA, the same three SNPs (rs10036748, rs6889239, and rs960709) show the highest

909    association values and are all in high LD. In EA rs6889239 has the best p-value and

910    rs10036748 yields the best p-value in HA.

911    (B) Analyzing the region with SKAT in a ∆MGW-weighted approach. In this region, for these

912    SNPs, including ∆MGW did not provide differential prioritization, rs6889239 remained the top

913    signal for EA and rs10036748 for HA.

914    (C) For each SNP in the 40 kb region, the posterior probabilities are plotted for the derived,

915    ∆MGW-weighted Bayesian analysis. The weighted Bayesian analysis did not alter the relative

916    signals observed in the single-SNP logistic regression analyses. In the EA data, rs6889239

917    yielded the largest posterior probability in EA and rs10036748 remained the top signal for HA.

918    (D) The ∆MGW is plotted for each genotyped SNP that passed quality control measures. SNPs

919    that were identified by the derived ∆MGW-weighted credible set are plotted in yellow. These

920    prioritized SNPs have comparatively low magnitudes of ∆MGW, indicating that the driving factor

921    of these SNP prioritizations stemmed from their SLE associations and not their magnitude of

922    ∆MGW.

923

924

925

926 **Primary Tables**

927 **Table 1. Summary statistics for the complete ΔMGW (Å) sample space.**

| 5th Nucleotide pairing[a] | N | Min. | Max. | Range | Median | Mean | Standard Deviation |
|---|---|---|---|---|---|---|---|
| A/C | 65,536 | 0.03 | 2.74 | 2.71 | 0.86 | 0.90 | 0.39 |
| A/G | 65,536 | 0.05 | 2.07 | 2.02 | 0.46 | 0.50 | 0.25 |
| A/T | 65,536 | 0.07 | 3.16 | 3.09 | 1.11 | 1.16 | 0.48 |
| C/G | 65,536 | 0.00 | 1.44 | 1.44 | 0.62 | 0.64 | 0.27 |
| C/T | 65,536 | 0.05 | 2.07 | 2.02 | 0.46 | 0.50 | 0.25 |
| G/T | 65,536 | 0.03 | 2.74 | 2.71 | 0.86 | 0.90 | 0.39 |
| All Possible | 393,216 | 0.00 | 3.16 | 3.16 | 0.67 | 0.77 | 0.42 |

928   [a]Pairings generated by 5th nucleotide in 9-mer sequence, all other nucleotides held constant.

929

930 **Table 2. Summary Statistics for ΔMGW (Å) across bi-allelic SNPs in dbSNP SNP150 dataset.**

| | SNP Category | N | Min. | Max. | Range | Median | Mean | Standard Deviation |
|---|---|---|---|---|---|---|---|---|
| | dbSNP SNP150 (bi-allelic) | 199,038,197 | 0.00 | 3.16 | 3.16 | 0.56 | 0.68 | 0.43 |
| Single-NCBI Function Subsets | coding-synonymous | 1,178,980 | 0.00 | 2.58 | 2.58 | 0.48 | 0.55 | 0.30 |
| | intron | 84,909,115 | 0.00 | 3.16 | 3.16 | 0.56 | 0.68 | 0.42 |
| | missense | 2,345,831 | 0.00 | 3.16 | 3.16 | 0.52 | 0.61 | 0.36 |
| | ncRNA | 499,593 | 0.00 | 3.16 | 3.16 | 0.54 | 0.63 | 0.38 |
| | near-gene-3 | 654,589 | 0.00 | 3.16 | 3.16 | 0.55 | 0.66 | 0.41 |
| | near-gene-5 | 2,487,192 | 0.00 | 3.16 | 3.16 | 0.54 | 0.65 | 0.41 |
| | nonsense | 66,275 | 0.00 | 3.16 | 3.16 | 0.55 | 0.65 | 0.37 |
| | splice-3 | 25,401 | 0.01 | 2.07 | 2.05 | 0.57 | 0.61 | 0.31 |
| | splice-5 | 28,983 | 0.00 | 2.74 | 2.74 | 0.57 | 0.65 | 0.31 |
| | stop-loss | 2,225 | 0.03 | 3.16 | 3.13 | 0.61 | 0.71 | 0.42 |
| | unknown | 99,004,130 | 0.00 | 3.16 | 3.16 | 0.57 | 0.69 | 0.43 |
| | untranslated-3 | 1,299,685 | 0.00 | 3.16 | 3.16 | 0.55 | 0.67 | 0.41 |
| | untranslated-5 | 181,208 | 0.00 | 3.16 | 3.16 | 0.50 | 0.58 | 0.33 |

931