1 **Title:**

2 Stability of DNA methylation and chromatin accessibility in structurally diverse maize genomes.

3

4 **Running title:**

5 Chromatin stability in diverse maize genomes

6

7 **Authors:**

8 Jaclyn M Noshay[1], Zhikai Liang[1], Peng Zhou[1], Peter A Crisp[2], Alexandre P Marand[3], Candice N

9 Hirsch[4], Robert J Schmitz[3], Nathan M Springer[1*]

10

11 **Author Affiliations:**

12 [1] Department of Plant and Microbial Biology, University of Minnesota, St, Paul, MN 55108

13 [2] School of Agriculture and Food Sciences, University of Queensland, Australia

14 [3] Department of Genetics, University of Georgia, Athens, GA, 30602

15 [4] Department of Agronomy and Plant Genetics, University of Minnesota, St. Paul, MN, 55108

16

17 *Corresponding author:

18 Nathan Springer

19 springer@umn.edu

20

21

22

**Abstract:**

Accessible chromatin and unmethylated DNA are associated with many genes and cis-regulatory elements. Attempts to understand natural variation for accessible chromatin regions (ACRs) and unmethylated regions (UMRs) often rely upon alignments to a single reference genome. This limits the ability to assess regions that are absent in the reference genome assembly and monitor how nearby structural variants influence variation in chromatin state. In this study, *de novo* genome assemblies for four maize inbreds (B73, Mo17, Oh43 and W22) are utilized to assess chromatin accessibility and DNA methylation patterns in a pan-genome context. The number of UMRs and ACRs that can be identified is more accurate when chromatin data is aligned to the matched genome rather than a single reference genome. While there are UMRs and ACRs present within genomic regions that are not shared between genotypes, these features are substantially enriched within shared regions, as determined by chromosomal alignments. Characterization of UMRs present within shared genomic regions reveals that most UMRs maintain the unmethylated state in other genotypes with only a small number being polymorphic between genotypes. However, the majority of UMRs between genotypes only exhibit partial overlaps suggesting that the boundaries between methylated and unmethylated DNA are dynamic. This instability is not solely due to sequence variation as these partially overlapping UMRs are frequently found within genomic regions that lack sequence variation. The ability to compare chromatin properties among individuals with structural variation enables pan-epigenome analyses to study the sources of variation for accessible chromatin and unmethylated DNA.

**Article summary:**

Regions of the genome that have accessible chromatin or unmethylated DNA are often associated with cis-regulatory elements. We assessed chromatin accessibility and DNA methylation in four structurally diverse maize genomes. There are accessible or unmethylated regions within the non-shared portions of the genomes but these features are depleted within these regions. Evaluating the dynamics of methylation and accessibility between genotypes reveals conservation of features, albeit with variable boundaries suggesting some instability of the precise edges of unmethylated regions.

**Introduction:**

The 2.1Gb maize B73 genome was first assembled in 2009 and contains ~80% repetitive sequence (Schnable *et al.* 2009). Unlike model species such as *Arabidopsis thaliana*, maize has transposable elements and highly methylated regions that are interspersed with genic regions of the genome (The Arabidopsis Genome Initiative 2000; Baucom *et al.* 2009; Springer and Schmitz 2017). One challenge in complex crop genomes such as maize is the identification of regulatory elements within genomes. There are opportunities to utilize both chromatin properties such as DNA methylation or chromatin accessibility to identify functional elements.

The maize genome is highly methylated and regions containing DNA methylation can be sub-classified based on the specific sequence context of the methylation. High levels of CG and CHG (H = A, C or T) methylation without CHH methylation are often found over transposable elements and other repetitive regions of the maize genome, while CG-only methylation is observed frequently within gene bodies (West *et al.* 2014; Niederhuth *et al.* 2016; Crisp *et al.* 2020). CHH methylation, which is largely the result of RNA-directed DNA methylation (RdDM), is found near highly expressed genes (Gent *et al.* 2013; Li *et al.* 2015a; Niederhuth *et al.* 2016). A small proportion of the maize genome lacks DNA methylation in any sequence context and these unmethylated regions (UMRs) likely reflect regions with potential roles in regulation of gene expression (Oka *et al.* 2017; Ricci *et al.* 2019; Hoefsloot and Stam 2020; Crisp *et al.* 2020).

Chromatin accessibility is another feature of chromatin that can be used to identify genomic regions with roles in regulation of transcription. In maize, ~1% of the genome contains accessible chromatin when profiled with a single tissue type (Rodgers-Melnick *et al.* 2016). Profiles of chromatin accessibility combined with other chromatin modifications have identified potential regulatory elements in the maize genome (Oka *et al.* 2017; Ricci *et al.* 2019). While chromatin accessibility is quite useful for identifying regulatory elements in a particular tissue, this property is highly dynamic with changes between tissue types or cells (Ricci *et al.* 2019; Marand *et al.* 2020a; Crisp *et al.* 2020). The vast majority of accessible chromatin occurs in unmethylated regions of the genome. However, there are additional unmethylated regions that do not exhibit chromatin accessibility. These likely reflect the fact that the unmethylated regions of the genome are quite stable in vegetative tissues while chromatin accessibility is highly tissue-specific (Schmitz *et al.* 2013; Kawakatsu *et al.* 2016; Marand *et al.* 2020b; Crisp *et al.*

88    2020). To date, the analysis of chromatin accessibility in maize has largely focused on the

89    accessible regions within the B73 genome.

90

91    The analysis of chromatin properties within the B73 reference genome has been useful for

92    functional annotation of the genome. However, there is also value in assessing natural variation

93    for the chromatin properties in different inbred lines of maize. While chromatin accessibility

94    studies have largely focused on B73, many studies have compared DNA methylation between

95    maize genotypes (Eichten *et al.* 2013; Regulski *et al.* 2013; Li *et al.* 2015b; Anderson *et al.*

96    2018; Xu *et al.* 2019, 2020). These studies have found many examples of DNA methylation

97    variation. Changes in DNA methylation can occur due to alterations in DNA sequence such as

98    transposon insertions (Noshay *et al.* 2019) or can occur in regions with no genetic changes

99    (Eichten *et al.* 2011). The ability to fully compare DNA methylation patterns among genotypes

100   and to investigate the role of structural variation has been limited due to reliance upon a single

101   reference genome for comparisons.

102

103   The genome content varies substantially among maize genotypes (Fu and Dooner 2002;

104   Springer *et al.* 2009; Swanson-Wagner *et al.* 2010; Anderson *et al.* 2019; Hufford *et al.* 2021).

105   The availability of multiple *de nov*o assembled reference genomes has enabled whole genome

106   comparisons of genome content (Hirsch *et al.* 2016; Springer *et al.* 2018; Sun *et al.* 2018;

107   Haberer *et al.* 2020; Hufford *et al.* 2021). Many of the sequences present in any one inbred are

108   not present at collinear regions in other genomes (Fu and Dooner 2002; Sun *et al.* 2018;

109   Haberer *et al.* 2020). This results in a pan-genome that contains more genes and transposons

110   than any individual maize inbred (Hirsch *et al.* 2014; Anderson *et al.* 2019; Hufford *et al.* 2021).

111   While it is quite clear that genome content differs substantially, it has been difficult to assess the

112   chromatin of the pan-genome due to technical difficulties in connecting the same sequence

113   regions between genotypes.

114

115   In this study, we generated DNA methylation and chromatin accessibility profiles from four

116   maize inbreds that each have *de novo* genome assemblies. UMRs and ACRs are identified for

117   each genotype based on alignment of the chromatin data to the B73v4 genome and the genome

118   from which it was generated. Chromosomal alignments were used to classify shared and

119   nonshared sequences between genomes. UMRs and ACRs are substantially depleted within

120   the non-shared portions of the genome. We assessed the stability of UMRs between genotypes

121   within the shared regions of the genome.  While the majority of UMRs in these regions have an

122    overlapping UMR in another genotype, the majority do not have identical boundaries.  These

123    UMRs with shifted boundaries account for a large portion of the differentially methylation regions

124    between two genotypes. The partially overlapping UMRs are not enriched for variable chromatin

125    accessibility or changes in expression of nearby genes, suggesting that differences in the

126    specific boundaries between methylated and unmethylated DNA are tolerated with little

127    functional impact.

128

129    **Results:**

130    ***Characterization of unmethylated DNA and accessible chromatin in four maize genomes***

131    DNA methylation (profiled using whole genome bisulfite sequencing - WGBS (Cokus *et al.* 2008;

132    Lister and Ecker 2009)), chromatin accessibility (profiled using Assay for Transposase

133    Accessible Chromatin-sequencing - ATAC-seq (Buenrostro *et al.* 2013)), and gene expression

134    (RNA-seq) data were generated for the same tissue sample from seedling leaf of four maize

135    inbreds (B73, Mo17, W22 and Oh43) (Table S1). For all genotypes the resulting datasets were

136    aligned to their own genome assembly and non-B73 genotypes were additionally aligned to the

137    B73v4 reference genome assembly.

138

139    The alignment rates for the WGBS datasets were substantially higher when mapped to their

140    respective genome assembly (~60%) compared to non-B73 samples mapped to the B73

141    reference genome assembly (~43%) (Table S1). The reduced mapping rate when aligning data

142    from non-B73 genotypes to the B73 genome assembly is likely due to polymorphisms and

143    structural variants present between inbreds. We focused on analysis of methylation

144    classifications based analysis of merged replicates, since the data from the two biological

145    replicates was highly correlated and the unmethylated regions identified within individual

146    samples were frequently (>97%) found in the merged sample (Table S3). The WGBS data was

147    used to classify the methylation state for each 100bp bin based on context-specific DNA

148    methylation (Figure S1A) as described previously (Crisp *et al.* 2020). Bins were classified as

149    unmethylated (<20% methylation in all contexts), CHH (CHH>15%), CG/CHG (>40% both CG

150    and CHG), CG only (>40% CG), missing data, missing sites or intermediate methylation (Figure

151    S1A). The majority (71-74%) of the maize genome is classified as methylated with most of this

152    exhibiting CG/CHG methylation and in rare cases CHH methylation (Figure S1A). A much

153    smaller proportion (6-7%) of the genome is classified as unmethylated (Figure S1A). In each

154    genome, roughly 15% of the bins are classified as missing data, likely due to an inability to align

155    WGBS reads uniquely to repetitive regions. However, the proportion of bins with missing data

156  was substantially larger when non-B73 WGBS data was aligned to the B73 genome (Figure

157  S1B).

158

159  The unmethylated 100bp bins were merged and filtered (Crisp *et al.* 2020) to identify

160  unmethylated regions (UMRs) (Table 1). UMRs were defined for each inbred based on

161  alignment to their respective genome assembly and when aligned to B73 (Figure 1A). The total

162  number of UMRs was similar across all four genotypes, although a greater number of UMRs

163  were defined when mapping WGBS reads to the sample-matched genome assembly. UMRs

164  were classified as genic, proximal (<2kb from nearest gene) and intergenic (>2kb from nearest

165  gene) in all four genotypes based on alignment of samples to their cognate reference genome

166  assembly, and were consistent across genotypes with ~50% of UMRs being observed in genic

167  regions and ~40% in intergenic regions  (Figure 1B).

168

169  Prior studies have found that unmethylated portions of the maize genome often contain cis-

170  regulatory regions (Oka *et al.* 2017; Ricci *et al.* 2019; Crisp *et al.* 2020). To determine the

171  concordance between UMRs and ACRs, we implemented ATAC-seq in the same four

172  genotypes. ACRs were identified in each individual sample as well as from merged biological

173  replicates (Table S2). We focused on analysis of the ACRs identified from the merged

174  replicates, since the data from the two biological replicates was highly correlated and the ACRs

175  identified within individual samples were frequently found in the merged sample (Table S3,

176  Figure S2). There are 21,232-24,309 ACRs present in each of the four genotypes (Table 1,

177  Figure 1C). Relative to UMRs, the ACRs are more enriched in gene-proximal regions of the

178  genome and depleted within intergenic regions, but >24% of the ACRs are found >2kb from the

179  nearest gene (Figure 1B). The vast majority of ACRs are found within UMRs in each of the four

180  genotypes (Figure 1D, S3). While the vast majority of ACRs occur within UMRs, there are many

181  UMRs without accessible chromatin (Figure 1D). This allows the classification of UMRs as

182  accessible UMRs (aUMRs) or inaccessible UMRs (iUMRs) based on whether they overlap an

183  ACR. The presence of an aUMR, which includes the presence of an accessible chromatin

184  region, is much more common within or near genes that are highly expressed, but is quite rare

185  for lowly expressed genes (Figure S3D). In contrast, iUMRs are present near genes with low

186  and high expression levels, but are depleted near silent genes (Figure S3E). While the aUMRs

187  represent an overlap between an unmethylated region and chromatin accessibility, the

188  boundaries of these regions are often not the same. The majority (97.3%) of cases represent a

189  larger unmethylated region in which the ACR only covers a portion of the UMR and the ACR is

190    often found in the center of the UMR (examples in Figure S4).  This suggests that the transition

191    from accessible to inaccessible chromatin and from unmethylated DNA to methylated DNA does

192    not occur at the same region.

193

194

195    *Classification of shared and non-shared genomic regions*

196    Previous studies have assessed natural variation in DNA methylation based on alignment to a

197    single reference genome (Regulski *et al.* 2013; Li *et al.* 2015b). However, when WGBS data

198    from non-B73 genotypes are aligned to the B73 genome, the proportion of regions with missing

199    data increases substantially (Figure S1B) and the methylation levels for genomic regions

200    missing in B73 are not assessed. The availability of multiple reference genomes provides the

201    opportunity to assess DNA methylation levels in the pan-genome that includes both shared

202    (syntenic) regions of the genome with or without allelic variation, as well as non-shared regions

203    that are present in one line and missing in another. The alignment of WGBS or ATAC-seq data

204    to their respective genome provides the advantage of more complete characterization of DNA

205    methylation and/or chromatin accessibility, but introduces complications for the direct

206    comparison of specific regions among genomes.

207

208    To address this complication in comparing regions across genomes, chromosomal alignments

209    were performed between the B73 genome and the other reference genomes to identify the

210    shared and non-shared genomic segments between any two genotypes (see Methods) (Figure

211    2A). The approach that was implemented employed relatively stringent criteria for identification

212    of shared regions. The regions classified as non-shared include both structural variants and

213    highly polymorphic regions as well as highly repetitive regions that could not be uniquely

214    mapped. Approximately 55% of the non-B73 genome sequences could be classified as syntenic

215    and mappable relative to B73, with the remaining 45% not aligning to the B73 genome due to

216    non-syntenic sequence or unmappable regions (Figure 2B). As a quality control measure, we

217    assessed the proportion of space classified as shared or non-shared within identity-by-state

218    (IBS) regions between genomes. The majority (94%) of IBS regions are classified as shared

219    between any two genomes (Table S4) and the regions that are not classified as shared within

220    IBS regions are highly enriched for repetitive sequences.

221

222    Our analysis of DNA methylation or chromatin accessibility is often focused on 100bp bins. To

223    directly compare the same coordinate space between genomes, we identified the 100bp bins

224     from the B73 genome that were shared across genotypes (Figure 2A, S4). In the comparisons

225     of B73 to the other three genomes, we find 41-48% of the B73 bins are non-shared, 37-42% of

226     bins have an exact match in shared regions, 12-14% mapped with >= 1 SNP, and an additional

227     4% mapped with >= 1 small (<20bp) indel between the two genotypes. Across all comparisons,

228     there are over ~800,000 100bp bins that are shared in all four genotypes (Figure 2C). There are

229     ~500,000 bins that are found only in B73 and another ~800,000 that are present in B73 and only

230     one or two of the other two genotypes (Figure 2C). The regions that are shared between

231     genotypes have fewer bins with missing data such that only 6.7% of the bins shared in all three

232     genotypes lack DNA methylation data compared to 28.4% of the bins that are only present in

233     B73. This likely reflects the fact that much of the non-shared sequence between genomes is

234     highly repetitive and recalcitrant to unique mapping. The identification of these shared bins

235     allowed us to calculate the methylation levels or ATAC-seq read depth for the specific

236     coordinates in a second genome that correspond to the B73 bins to allow direct comparisons of

237     chromatin properties between genomes using epigenomic data aligned to its own reference

238     genome.

239

240     ***UMRs and ACRs are depleted in non-shared portions of the genome***

241     We initially focused on the chromatin properties of the non-shared portions of the genome to

242     assess the frequency of UMRs or ACRs within the dispensable portion of the genome compared

243     to the shared portions. While over 10% of the shared genomic regions are annotated as genic

244     less than 4.8% of the non-shared regions are annotated as genic reflecting a depletion of genes

245     and enrichment of intergenic and TE sequence. The analysis of the *bronze1 (bz1)* locus on

246     chromosome 9 illustrates these trends of shared space in genic regions and large non-shared

247     blocks between genes, as previously described (Fu and Dooner 2002; Wang and Dooner 2006)

248     (Figure 3A). In the *bz1* region, very few UMRs or ACRs are found within the non-shared regions

249     (Figure 3B). We proceeded to perform a genome-wide assessment of the proportion of UMRs

250     within shared and non-shared regions of the genome. While UMRs account for 6% of the entire

251     B73 genome, only ~2% of the non-shared genomic regions are classified as UMRs compared to

252     ~12% of the shared genomic regions (Figure 3C).  A similar analysis of the genome-wide

253     distribution of ACRs reveals that accessible chromatin is even more enriched within genomic

254     regions that are shared among all four genotypes (Figure 3C). ACRs account for 1.2% of the

255     shared genomic space but only 0.1% of the non-shared genomic regions (Figure 3C).  Both

256     ACRs and UMRs tend to be enriched near genes and the non-shared genomic regions are

257     relatively gene-poor.  However, this depletion of genes is not the only explanation for the paucity

258   of ACRs and UMRs in the non-shared genomic regions.  Over 80% of the genes in the shared

259   space contain a UMR, while only 17% of the genes located in non-shared regions contain

260   UMRs.  Prior studies have found that non-shared genes are less likely to be expressed (Hirsch

261   *et al.* 2016; Sun *et al.* 2018; Anderson *et al.* 2019; Haberer *et al.* 2020) and the depletion of

262   UMRs within or near these genes further suggests that many of these "genes" lack the

263   chromatin properties associated with expression. These analyses suggest that pan-genome

264   assessment of UMRs and ACRs will provide a more complete identification of UMRs/ACRs but

265   that there are a limited number of novel UMRs or ACRs in non-shared space in maize. The

266   subsequent analysis will focus on the UMRs and ACRs that are present within shared regions of

267   any two maize genomes.

268

269   ***Comparisons of UMRs and ACRs in the shared space of maize genomes***

270   We proceeded to focus on the UMRs and ACRs that are present within shared regions between

271   maize genomes. The analyses were primarily focused on UMRs as these encompass the vast

272   majority of ACRs (Figure 1D) and we could monitor stability for the UMRs with an ACR (aUMRs)

273   compared to the UMRs without an ACR (iUMRs). The B73 UMRs were compared to each of the

274   other genomes and classified based on whether they are present in shared/non-shared regions

275   and then whether the region has DNA methylation data available for both genotypes.  For the

276   ~90% of B73 UMRs that have defined methylation states and are present in a shared region we

277   could classify whether there is an overlapping UMR in the other genotype or whether the UMR

278   is polymorphic such that it is classified as methylated in the other genotype (Figure 4A). Most

279   UMRs that are present in shared space overlap a UMR in the other genotype while a small set

280   are polymorphic (Figure 4A).  The overlapping UMRs can be classified as identical if the

281   boundaries of the UMR are the same in both genotypes (example in Figure 4B).  Alternatively,

282   an overlapping UMR could represent a partial overlap such that one genotype has a larger

283   region than the other or both edges are shifted (examples in Figure 4B). The UMRs with partial

284   overlap account for the majority of the overlapping UMRs between two genotypes (Figure 4A).

285

286   B73 UMRs can be subdivided into aUMRs and iUMRs based on the presence, or absence, of

287   an ACR within the unmethylated region. We compared the distribution of classifications for the

288   aUMRs and iUMRs for the presences of identical, partially overlapping, or polymorphic UMRs in

289   the other genotypes (Figure 4C). The B73 aUMRs have fewer examples of polymorphic UMRs

290   as well as fewer examples within non-shared genomic regions. However, this is largely due to a

291   larger proportion of overlapping UMRs that are partially overlapping rather than more examples

292    of identical UMRs (Figure 4C). These analyses suggest that while any two genomes often have

293    UMRs in similar regions the exact coordinates of the UMRs are often distinct.

294

295    The B73 aUMRs and iUMRs were also assessed for the potential changes to either methylation

296    or accessibility between genotypes (Figure 4D). The majority (~71.2%) of the B73 aUMRs were

297    maintained as aUMRs in the other genotypes.  However, there are also a subset of the B73

298    aUMRs that lose either the unmethylated state (~14.8%) or chromatin accessibility (~11.1%) in

299    the other genotype.  The remaining 2.9% are not classified as either ACR or UMR for the same

300    region in the other genome.  The B73 iUMRs often (~73.7%) are unmethylated and inaccessible

301    in the other genotypes (Figure 4D). There are also many (~24.5%) examples of B73 iUMRs that

302    are methylated in the other genotype.  The proportion of shifts from unmethylated to methylated

303    states are much higher for the iUMRs than the aUMRs.  Very few (~1.6%) of the B73 iUMRs

304    exhibit accessibility in the other genotypes (Figure 4D).

305

306    ***Unique properties of regions with methylation changes in various methylation contexts***

307    While the polymorphic B73 UMRs that are methylated in another genotype only account for a

308    small set of all UMRs these may represent important functional differences between genotypes.

309    The polymorphic UMRs can be subdivided based on the prominent class of methylation in the

310    other genotype (Figure 4A, 5A). Each of these classes of methylation gains likely reflect distinct

311    mechanisms and chromatin types. The types of methylation observed in these regions do not

312    reflect the genome-wide proportions of methylation types (Figure S1). The proportions that are

313    classified as CG only or CHH are higher than observed genome wide (Figure S1, 5B). While

314    CG/CHG regions are depleted, although there are still many examples of CG/CHG at these

315    regions of variable methylation (Figure 5B).

316

317    The presence or absence of ACRs in both genotypes was assessed for the polymorphic UMRs

318    relative to overlapping UMRs (Figure 5C). While both identical UMRs and partially overlapping

319    UMRs show virtually identical proportions with shared ACRs or polymorphic ACRs, the

320    polymorphic UMRs have very few stable ACRs (Figure 5C).  This is expected as there are very

321    few examples of accessible regions within methylated DNA. As expected, there are very few

322    examples of ACRs only in the methylated genotype. The proportion of the polymorphic UMRs

323    that are classified as having an ACR only in B73 but not the methylated genotype is quite

324    variable. Polymorphic UMRs with CHH methylation in the other genotype are more likely to have

325    an ACR in B73 than polymorphic UMRs with CG only methylation in the other genotype (Figure

326    5C).  This could reflect the fact that CHH methylation is often found in regions immediately

327    upstream or downstream of genes in the maize genome (Gent *et al.* 2013; Li *et al.* 2015a) and

328    that these regions often have ACRs.  In contrast the CG-only methylation often occurs within

329    gene bodies, where ACRs are less common than at the edges of genes or promoter regions.

330

331    We proceeded to assess variable gene expression of genes near overlapping or polymorphic

332    UMRs using RNAseq data from the same tissue used to monitor accessibility and DNA

333    methylation.  Genes with an overlapping or polymorphic UMR within 200bp (upstream or

334    downstream) of the transcription start site (TSS) were identified and classified as being

335    differentially expressed (DE), expressed in both genotypes but not DE, or silent (FPKM < 1 in

336    both genotypes). Genes that have identical or partially overlapping UMRs near the TSS exhibit

337    nearly identical proportions of genes in these categories and have similar proportions of genes

338    that are higher expressed in B73 or the other genotype (Figure 5D). Polymorphic UMRs that

339    gain CG-only methylation in the other genotype have fewer examples of genes with higher

340    expression in B73.  This suggests that the presence of CG only methylation is rarely associated

341    with reduced gene expression.  In contrast, genes with gains of CG/CHG or CHH methylation

342    near the TSS are enriched for genes that are higher expressed in B73. While there is an

343    enrichment for DE expression with the unmethylated genotype being higher expression for

344    polymorphic UMRs with CG/CHG or CHH methylation gains, it is worth noting that there are

345    also many examples of genes near these types of polymorphic UMRs that are not differentially

346    expressed.  This suggests that the gain of CG/CHG methylation or CHH methylation in regions

347    surrounding the TSS can be associated with altered expression in some cases, but that other

348    genes can tolerate variable methylation without a significant change in expression.

349

350    ***Partially overlapping UMRs contribute substantially to differentially methylated regions***

351    The analysis of natural variation for DNA methylation is often focused on identification of

352    differentially methylated regions (DMRs) between genotypes.  In this study, we elected to focus

353    on the conservation / variation of unmethylated regions as these regions have evidence for

354    functional relevance in crop genomes. However, the observation that many of these regions

355    only have partial overlap suggests that many DMRs might be the result of a shift in the

356    boundary between methylated and unmethylated DNA rather than a complete regional gain/loss

357    of methylation (Figure 6A).  The 100bp bins were used to identify DMRs between the

358    genotypes. There are 116,000-158,000 100bp bins that are classified as differentially

359    methylated with hypomethylation in B73 relative to the other genotype. We assessed how many

360    of these DMRs are due to completely polymorphic UMRs compared to partial UMRs with

361    different boundaries between methylated and unmethylated DNA (Figure 6B).  The polymorphic

362    UMRs account for 2.5-3.3% of all differentially methylated bins depending on which genotypes

363    are being compared.  A larger proportion (51.5-53.5%) of the differentially methylated bins are

364    due to partially overlapping UMRs.  The remaining differentially methylated bins occur in regions

365    too small to be classified as UMRs (unmethylated regions <300bp) or represent single bin

366    differences in larger UMRs.  This analysis suggests that many of the DMRs are due to shifting

367    boundaries between methylated and unmethylated DNA rather than a complete gain or loss of

368    methylation in a region.

369

370    These observations suggest that the specific boundary between methylated and unmethylated

371    DNA can be variable between genotypes.  This could be due to sequence changes at or near

372    the edges of these regions or could arise due to stochastic variation with no sequence change.

373    To address this question we assessed the proportion of identical or partially overlapping UMRs

374    within large (>1Mb) blocks of sequence that is IBS.  In total there was 112.7Mb of IBS sequence

375    blocks that could be assessed and these are large blocks of sequence that are essentially

376    devoid of SNPs or structural variants.  Within these regions we find a depletion for polymorphic

377    UMRs. While 5.3% of all UMRs are classified as polymorphic we find only 2.8% of UMRs that

378    are classified as polymorphic in these regions suggesting that fully polymorphic UMRs are

379    depleted in the absence of sequence variation (Figure 6C). The IBS regions have a higher

380    proportion of UMRs with identical boundaries in the two genotypes.  However, there are still a

381    large number of UMRs with shifted boundaries (49.5%) suggesting that the boundaries between

382    methylated and unmethylated DNA can shift even without nearby sequence variation.

383

384    **Discussion:**

385    *Zea mays,* unlike many other model organisms, has a large genome containing 80% repetitive

386    sequence and high levels of DNA methylation interspersed with functional genic and regulatory

387    regions (Schnable *et al.* 2009; Jiao *et al.* 2017).  Examination of genome structure across inbred

388    lines have identified extensive polymorphism in both genic and repeat regions of the maize

389    genome (Chia *et al.* 2012; Hirsch *et al.* 2014; Springer *et al.* 2016; Darracq *et al.* 2018;

390    Anderson *et al.* 2019; Hufford *et al.* 2021).  Prior analyses of natural variation of chromatin in

391    maize have been based on epigenome profiling data aligned to a single reference genome (Li *et*

392    *al.* 2015b; Xu *et al.* 2020).  While a single reference genome provides insight into variation in

393    conserved genomic regions, it does not contain the full set of sequences present in the lines

394    being compared, resulting in biases in the ability to compare chromatin properties. The

395    availability of multiple *de novo* genome assemblies allows for a more complete discovery of

396    regions with specific chromatin properties, such as UMRs or ACRs. In this study, we profiled

397    genome-wide DNA methylation, based on alignments of data to the corresponding genome

398    assembly, to identify the ~6% of each genome that exhibits an unmethylated state and the ~1%

399    that is accessible chromatin. A pan-genomic analysis of UMRs and ACRs reveals the frequency

400    of these features within both shared and nonshared genomic regions. Within the shared

401    sequence regions it is possible to assess the stability of the unmethylated and accessible

402    chromatin portions of the genome.

403

404    ***Pan-genome analyses reveal enrichment of unmethylated regions within shared***

405    ***sequence***

406    Whole genome alignments between B73 and Mo17, W22, and Oh43 allowed for the

407    identification of both shared and nonshared sequences. In a comparison of any two genomes,

408    the sequence unique to each genome is primarily composed of highly repetitive sequences with

409    extensive DNA methylation and is found to be depleted for genes (Chia *et al.* 2012; Springer *et*

410    *al.* 2016; Hirsch *et al.* 2016; Darracq *et al.* 2018; Anderson *et al.* 2019; Hufford *et al.* 2021).  The

411    proportion of the nonshared genome that is classified as UMR or ACR is 6-12 fold lower than

412    the proportion of the shared genome classified as UMR or ACR.  This reduction in UMRs and

413    ACRs is not simply due to the reduced gene content in nonshared space.  Most (80%) of genes

414    in the shared space are associated with a UMR, while only 17% of genes in non-shared space

415    have a UMR. This is not unexpected as prior studies of presence-absence variation (PAV)

416    genes have found that most of these genes that vary between genotypes are not expressed

417    even when they are present (Swanson-Wagner *et al.* 2010).  A recent analysis of 26 maize

418    genomes that used a slightly different approach to classify unmethylated and CG-only regions

419    reported similar findings (Hufford *et al.* 2021). More UMRs are identified by aligning chromatin

420    data to the proper genome but the proportion of UMRs or ACRs in this nonshared space is

421    much lower than in the shared regions. These analyses suggest that pan-genomic analyses can

422    identify novel UMRs or ACRs but that these are relatively rare in the sequences that exhibit

423    large scale structural variation. However, it is worth noting that the UMRs or ACRs that are

424    present near the genes in the non-shared space can be an effective tool for identifying genes

425    with potential expression (Sartor *et al.* 2019; Crisp *et al.* 2020).  Given that many of the genes

426    within these regions are likely pseudogenes generated by transposition of genes or gene

427    fragments that can be difficult to annotate just based on sequence, the use of chromatin data

428    can help to identify genes with potential function in these regions.

429

430    ***Characterization of relative dynamics of accessibility and methylation***

431    We were interested in studying the relative dynamics of both DNA methylation and chromatin

432    accessibility among genotypes.  Prior studies have found that the majority of accessible regions

433    have little or no methylation (Ricci *et al.* 2019) but that there are also many unmethylated

434    regions that lack accessibility (Crisp *et al.* 2020). The analysis of UMRs that are present within

435    shared sequence regions can be used to understand how often there is variation in only

436    accessibility as opposed to coordinate changes in both DNA methylation and accessibility. The

437    accessible UMRS (aUMRs) tend to be relatively stable in other genotypes with both accessibility

438    and lack of DNA methylation for an overlapping region in other haplotypes.  This is consistent

439    with the concept that these regions may be important for proper regulation of gene expression

440    and therefore changes in these chromatin properties could be associated with functional

441    differences. The inaccessible UMRs (iUMRs) were often inaccessible and unmethylated in both

442    genotypes but there were a large number of these that exhibit polymorphic DNA methylation

443    status such that they exhibit high levels of DNA methylation in the other genotypes.  Only a

444    small proportion of these UMRs exhibit a consistently unmethylated state in both genotypes with

445    accessible chromatin in only one of the two genotypes. These likely include some examples of

446    false negatives due to relatively stringent criteria for calling an ACR.  In these cases, an ACR

447    may be present in both genotypes but only identified as significant for one genotype. However,

448    these cases of variable chromatin accessibility also include examples with clear support for

449    chromatin accessibility in one genotype but no evidence for chromatin accessibility in the other

450    genotype.  These are interesting as they potentially reflect differences in transcription factor

451    occupancy for regions that are stably unmethylated in both genotypes. It is possible that these

452    may reflect differences in tissue-specific expression patterns of some maize genes.  In leaf

453    tissue there may be differential chromatin accessibility, but it is possible that the genotype

454    without chromatin accessibility in leaf tissue still becomes accessible in some other tissue that

455    exhibits expression.  Alternatively, minor sequence changes at transcription factor binding sites

456    may result in loss of chromatin accessibility even though the region is unmethylated in both

457    genotypes.

458

459    ***Stability and instability of UMRs between genotypes***

460    A subset of the shared sequence UMRs do not maintain their unmethylated state across

461    genotypes and instead have high levels of methylation in at least one of the other three

462    genotypes.  The presence of methylation variation in the shared sequence regions allowed for

463    characterization of attributes associated with chromatin state instability. Prior studies have

464    suggested that structural variants, especially transposable element polymorphisms, can be

465    associated with changes in DNA methylation for nearby sequences (Eichten *et al.* 2012;

466    Schmitz *et al.* 2013).  When analyzing DNA methylation based on a single reference genome it

467    can be difficult to incorporate information about structural variants and to map reads near the

468    junctions of these variants.  Using alignments to each reference genome and then comparing

469    coordinates of syntenic 100bp tiles allowed us to monitor changes in DNA methylation between

470    genotypes, even in regions near structural variants.  The polymorphic UMRs that represent a full

471    shift of an unmethylated region in one genotype to methylation in the other genotype are

472    depleted in regions devoid of structural variants.  Within large blocks of IBS 2.8% of the UMRs

473    are polymorphic.  In contrast, over 5.3% of all UMRs are classified as polymorphic.  This

474    indicates that changes in methylation state can occur in the absence of nearby structural

475    variants but that the rate is substantially higher in regions with sequence variation.

476

477    In this study, we focused on the conservation and variation for UMRs or ACRs between

478    genotypes.  These are relatively large (at least 300bp based on the criteria used for discovery)

479    regions that lack DNA methylation.  We focused on these regions due to prior evidence for

480    functional enrichment of these regions (Oka *et al.*; Ricci *et al.* 2019). We note that most of the

481    UMRs in one genotype have an overlapping UMR in another genotype.  This suggested stability

482    of these chromatin patterns among genotypes.  However, closer inspection revealed that the

483    majority of these overlapping UMRs have different boundaries in the two genotypes.  These

484    include examples in which one UMR is entirely within the other as well as examples that have

485    partial overhangs in both genotypes.  The partially overlapping UMRs seem to have very similar

486    genomic distributions and overlap with ACRs or altered gene expression in similar proportion to

487    those for identical conserved UMRs.  This suggests that these shifts in the boundary between

488    methylated and unmethylated DNA do not have functional impact in most cases. This may

489    suggest that the presence of a UMR is more defined by sequences in the middle of the

490    unmethylated region rather than particular sequences at the edges that define the extent of

491    methylation.

492

493     The observation of many partially overlapping UMRs suggested that these shifts in the

494     boundary between methylated and unmethylated DNA could account for many examples of

495     differential methylation between genotypes.  Conceptually, it is tempting to think that most

496     differentially methylated regions result from a local gain or loss of a patch of DNA methylation.

497     However, our analyses suggest that many of the differentially methylated 100bp tiles actually

498     arise due to changes in the boundaries between UMRs in different genotypes. Further studies

499     will be necessary to determine if these differences in methylation boundaries represent a

500     continuum such that each genotype has a slightly different boundary or if there are preferred

501     epi-haplotypes.

502

503     **Methods:**

504     ***Reference Genomes:***

505     Whole genome assemblies for four maize inbred lines, B73 (Jiao *et al.* 2016), W22 (Springer *et*

506     *al.* 2018), Mo17 (Sun *et al.* 2018), and Oh43 (Hufford *et al.* 2021) were used for genome-wide

507     analyses. All analyses were performed on assemblies of chromosomes 1-10 while all unplaced

508     scaffolds were disregarded due to the inability to compare these regions across genotypes.

509     Filtered gene and structural TE annotations (Stitzer *et al.*; Anderson *et al.* 2019) were used.

510

511     ***Sample Collection:***

512     Maize B73, W22, Mo17 and Oh43 plants were grown under 16 h/8 h 30°C /20°C day/night for

513     13 days in the growth chamber of University of Minnesota.  DNA was extracted from leaves of

514     two-week old V2 plants using the DNeasy Plant Mini kit (Qiagen). Four or five biological

515     replicates consisting of a pool of tissue from 4 plants were collected for each genotype.  Two of

516     these biological replicates were sampled for profiling of DNA methylation and chromatin

517     accessibility while all biological replicates were used for RNAseq.

518

519     ***WGBS protocol:***

520     Two technical replicates of each genotype (B73, Mo17, W22, and Oh43) were generated.  1ug

521     of DNA in 50ug of water was sheared using an Ultrasonicator to approximately 200-350bp

522     fragments. 20ul of sheared DNA was then bisulfite converted using the EX DNA Methylation-

523     Lightning Kit (Zymo Research) as per the manufacturer's instructions and eluted in a final

524     volume of 15ul. Then 7.5ul of the fragmented bisulfite-converted sample was used as input for

525     library preparation using the ACCEL-NGS Methyl-Seq DNA Library Kit (SWIFT Biosciences).

526     Library preparation was performed as per the manufacturer's instructions. The indexing PCR

527     was performed for 5 cycles. Libraries were then pooled and sequenced on a NovaSeq 6000 in

528     high output mode 125bp paired end reads over a single lane at the University of Minnesota

529     Genomics Center. WGBS data generated in this study is deposited at NCBI SRA and available

530     under accession.

531

532     Trim_glore(Martin 2011) was used to trim adapter sequences and read quality was assessed with

533     the default parameters in paired-end read mode plus a hard clip of 20bp on each read due to

534     SWIFT protocol specifications. Reads that passed quality control were aligned to their

535     corresponding genome assemblies. Alignments were conducted using BSMAP-2.90(Xi and Li

536     2009), allowing only unique hits with up to 5 mismatches and a quality threshold of 20 (-v 5 -q

537     20). Duplicate reads were detected and removed using picard-tools-1.102 ("Picard") and

538     SAMtools(Li *et al.* 2009). Conversion rate was determined using the reads mapped to the

539     unmethylated chloroplast genome. The resulting alignment file, merged for all samples with the

540     same tissue and genotype, was then used to determine methylation level for each cytosine using

541     BSMAP tools.

542

543     ***Methylation data summary:***

544     Methylation levels were summarized using the bsmap methratio.py script to group by context (CG,

545     CHG, CHH). The number of cytosines in every 100bp bin of the genome was determined and the

546     proportion of cytosines defined as methylated was calculated. Coverage was calculated as CT /

547     # of sites for each context. Methylation domain was classified for each 100bp bin based on the

548     protocol described in Crisp et al. (Crisp *et al.* 2020) with criteria defined as a minimum site count

549     of 2 and coverage of 3. UMRs were defined by grouping adjacent unmethylated bins.

550

551     ***ATAC-seq protocol and ACR classification:***

552     ATAC-seq libraries were generated as described in Lu et al (Lu *et al.* 2017). Two technical

553     replicates of each genotype (B73, Mo17, W22, and Oh43) were generated from the same samples

554     as those used for WGBS data generation. Raw reads per sample were preprocessed with

555     Trim_gloare. Trimmed reads were aligned to the *Zea mays* B73v4 genome and the genome

556     assembly specific to each sample using Bowtie v1.2.3 with the following parameters: "bowtie -X

557     1000 -m 1 -v 2 --best –strata". Aligned reads were converted to bam files and sorted using

558     SAMtools v1.9. Clonal duplicates were removed using Picard MarkDuplicates v2.23.3

559     (http://broadinstitute.github.io/picard/). Input data of maize B73 was retrieved from a previous

560     publication and processed to obtain bam files with clonal duplicates removed. MACS2 was

561 employed to call initial ACRs with Input data as control (-c) and sample data as treatment (-t)
562 using the following parameter "-g 2.1e9 --keep-dup all --nomodel --extsize 147". The post-
563 processing followed the same procedure as a prior publication (Ricci *et al.* 2019) to produce high-
564 confident ACRs. Specifically, 1) Initial ACRs were split into 50 bp windows with 25 bp steps; 2)
565 the Tn5 integration frequency in each window was calculated and normalized to the average
566 frequency in the total genome; 3) windows with the normalized frequency greater than 25 were
567 merged together allowing 150 bp gaps; 4) only merged regions greater than 50 bp were retained;
568 5) the mitochondrial or chloroplast genome from NCBI Organelle Genome Resources were
569 removed using blast again sequences within merged ACR regions. The sites within ACRs that
570 had the highest Tn5 integration frequency were defined as summits.

571

572 ### *RNA-seq protocol:*
573 RNA-seq data were generated in 150bp paired-end mode using NovaSeq 6000. B73, W22 and
574 Mo17 reads were retrieved from the NCBI SRA accession PRJNA657262 (Liang et al. 2021) and
575 Oh43 reads were deposited into NCBI SRA accession PRJNA692023. All of the raw reads were
576 preprocessed using Trim_galore and aligned against the B73 AGPv4 reference genome using
577 HISAT2 v2.1.0 (Kim *et al.* 2015). Gene annotations and disjoined TE annotations were used as
578 described above. Gene exon regions were subtracted from TE regions and then appended to
579 original TE annotation to remove ambiguous mapping between genes and TEs. Reads per gene
580 or TE was determined using HTSeq-count v0.11.2 (Anders *et al.* 2015) and raw count data was
581 input into DESeq2 (Love *et al.* 2014) to identify differentially expressed genes or TE elements.

582

583 The mean value for each feature (gene or TE) was calculated from 4 or 5 replicates.  Any feature
584 with a mean value greater than 1 was considered "expressed". UMRs were associated with genes
585 and TEs based on location relative to the feature. B73 UMRs which overlapped the annotated
586 sequence coordinates within the genome being assessed were classified as "genic" or "TE".
587 Those not overlapping a gene but within 2kb of the gene start or end were classified as "proximal".

588

589 ### *Cross-genotype mapping:*
590 Genome sequence from Mo17, W22 and Oh43 was first aligned to the B73 reference (Jiao *et al.*
591 2017) using minimap2 (Li 2018). The resulting alignments were merged and cleaned (removing
592 overlapping alignment blocks and alignment blocks containing assembly gaps) using in-house
593 perl scripts. BLAT Chain/Net tools were then used to create a single coverage best alignment net
594 between the query genome (one of Mo17, W22 and Oh43) and the target genome (B73). Finally,

595    a genome-wide synteny chain file was built for each genotype (against HM101), enabling

596    downstream analyses such as variant detection and 100-bp tile liftover. Alignment pipeline and

597    scripts are available on Github (https://github.com/baudisgroup/segment-liftover). Sequence was

598    extracted for all 100bp bins in the B73 genome and aligned to Mo17, W22, and Oh43.  Each bin

599    was determined to be unmappable or mappable.  Mappable bins were assigned coordinates in

600    the non-B73 genome.  The number of single nucleotide polymorphisms and insertion/deletions

601    for each bin was calculated.  Across all genotypes, only 4% of bins were found to have >= 1

602    insertion/deletion and 13% contained >=1 single nucleotide polymorphism.  Bins with no more

603    than 4 insertion/deletions of 20bp in size were kept for analyses of shared space.  Each 100bp

604    bin in B73 was designated as unmapped or provided matching sequence coordinates in each of

605    the 3 other genotypes (Mo17, W22, Oh43).

606

607    ***Differentially methylated tiles (DMTs):***

608    WGBS data aligned to the respective genome and summarized in the B73-based 100bp

609    coordinate system was used.  Tiles were subset to those with sequence mappability and

610    coverage in both genotypes for each pairwise comparison.  DMTs were defined by a difference

611    of 40% with at least one genotype having <10% and >40% methylation for CG and CHG

612    contexts.  CHH DMTs were defined by one genotype with <5% and >25% methylation in the

613    100bp tile.  DMTs in each context were determined for Mo17, W22, and Oh43 compared to B73.

614

615    ***Classification of UMR variability:***

616    B73 UMRs that were mappable to sequence in another genotype were further defined by

617    methylation state in the corresponding genome.  All 100bp bins within a defined UMR were

618    assessed for the matching sequence coordinates in Mo17, W22, and Oh43.  For each UMR, the

619    proportion of bins classified as methylated (including CG, CG/CHG, and CHH methylation

620    domains) was calculated.  UMRs with >50% of the bins being methylated were defined as

621    "polymorphic UMRs" for the difference in methylation state from unmethylated in B73 to

622    methylated in the non-B73 genotype.  All other UMRs, showing an unmethylated state in both

623    B73 and the non-B73 genotype assessed, were defined as "overlapping UMRs".

624

625    B73 UMRs that are methylated in another genotype (polymorphic UMRs) were further classified

626    by the type of methylation observed in the non-B73 genotype.  The polymorphic UMRs were

627    summarized by domain.  The proportion of 100bp bins with a methylated domain, within the

628    defined B73 UMR, for each methylation context was determined.  Any UMR that had >50% of

629    its methylated bins classified as a specific methylation context was declared to be variable in

630    that context.  Classification was determined first by CHH methylation, followed by CG/CHG

631    methylation and lastly CG only methylation.  Variable methylation type was defined individually

632    for each genome based on the sequence coordinates of the B73 UMR.

633

634    B73 UMRs that are unmethylated in another genotype (overlapping UMRs) were further

635    classified by the coordinates of the defined UMR between B73 and the non-B73 genotype.  The

636    UMRs, defined by alignment of WGBS data to the B73 reference genome, were determined and

637    their coordinates were assessed.  Pairwise comparisons were done between B73 and non-B73

638    genotypes.  B73 UMRs that had identical 100bp bin boundaries for the defined UMR were

639    classified as identical UMRs.  B73 UMRs that had variable boundaries were classified as partial

640    UMRs (the coordinates of the smaller UMR were maintained within the larger UMR coordinates

641    or the coordinates are shifted and have uniquely defined unmethylated bins in each genotype).

642

643    ***Classification of ACR variability:***

644    Every B73 UMR was classified based on the accessibility of that shared sequence region within

645    B73, Mo17, W22, and Oh43.  All UMRs in B73 were defined as accessible (aUMR) or

646    inaccessible (iUMR) based on its overlap with an accessible chromatin region in the B73

647    sample.  For B73 aUMRs, the presence of an accessible region in the non-B73 genotypes was

648    determined.  The B73-based coordinates of the UMR in the corresponding genome were used

649    to identify overlap with the ACRs defined in that genome.  UMRs that overlap both an ACR in

650    B73 and non-B73 genome were defined as stable ACRs.  If the aUMR in B73 lacked

651    accessibility in the non-B73 genome it was defined as B73-only ACR.  Alternatively, if a UMR

652    was inaccessible in B73 it could never be found accessible or show accessibility in the other

653    genotype.  If the iUMR lacked accessibility in the non-B73 genome, it was determined to have

654    no ACR.  If the sequence of the iUMR overlapped a defined ACR in the other genome, it was

655    defined as a non-B73 ACR such that it was inaccessible in the B73 UMR but accessible in the

656    shared sequence of Mo17, W22, or Oh43.  The ACRs which were defined as either B73-only or

657    nonB73-only were verified by assessing the 100bp cpm values within that region across the two

658    genotypes.

659

660    **Data Availability Statement:**

661    Accessible chromatin data (ATAC-seq) generated for this study is available at NCBI short read

662    archive under accession number PRJNA709664.  In this study we also utilize previously

663   published RNA-seq datasets that are available under accession numbers PRJNA657262 and

664   PRJNA692023 and whole genome bisulfite datasets that are available under accession number

665   PRJNA657677.

666

675

676   **<u>References</u>**

677   Anderson, S. N., M. C. Stitzer, A. B. Brohammer, P. Zhou, J. M. Noshay *et al.*, 2019

678       Transposable elements contribute to dynamic genome content in maize.

679   Anderson, S. N., G. Zynda, J. Song, Z. Han, M. Vaughn *et al.*, 2018 Subtle Perturbations of the

680       Maize Methylome Reveal Genes and Transposons Silenced by Chromomethylase or RNA-

681       Directed DNA Methylation Pathways. G3 .

682   Anders, S., P. T. Pyl, and W. Huber, 2015 HTSeq--a Python framework to work with high-

683       throughput sequencing data. Bioinformatics 31: 166–169.

684   Baucom, R. S., J. C. Estill, C. Chaparro, N. Upshaw, A. Jogi *et al.*, 2009 Exceptional diversity,

685       non-random distribution, and rapid evolution of retroelements in the B73 maize genome.

686       PLoS Genet. 5: e1000732.

687   Buenrostro, J. D., P. G. Giresi, L. C. Zaba, H. Y. Chang, and W. J. Greenleaf, 2013

688       Transposition of native chromatin for fast and sensitive epigenomic profiling of open

689       chromatin, DNA-binding proteins and nucleosome position. Nat. Methods 10: 1213–1218.

690   Chia, J. M., C. Song, P. J. Bradbury, D. Costich, N. de Leon *et al.*, 2012 Maize HapMap2

691       identifies extant variation from a genome in flux. Nat. Genet. 44: 803–807.

692 Cokus, S. J., S. Feng, X. Zhang, Z. Chen, B. Merriman *et al.*, 2008 Shotgun bisulphite

693  sequencing of the Arabidopsis genome reveals DNA methylation patterning. Nature 452:

694  215–219.

695 Crisp, P. A., A. P. Marand, J. M. Noshay, P. Zhou, Z. Lu *et al.*, 2020 Stable unmethylated DNA

696  demarcates expressed genes and their cis-regulatory space in plant genomes. Proc. Natl.

697  Acad. Sci. U. S. A.

698 Darracq, A., C. Vitte, S. Nicolas, J. Duarte, J.-P. Pichon *et al.*, 2018 Sequence analysis of

699  European maize inbred line F2 provides new insights into molecular and chromosomal

700  characteristics of presence/absence variants. BMC Genomics 19: 119.

701 Eichten, S. R., R. Briskine, J. Song, Q. Li, R. Swanson-Wagner *et al.*, 2013 Epigenetic and

702  genetic influences on DNA methylation variation in maize populations. Plant Cell 25: 2783–

703  2797.

704 Eichten, S. R., N. A. Ellis, I. Makarevitch, C. T. Yeh, J. I. Gent *et al.*, 2012 Spreading of

705  heterochromatin is limited to specific families of maize retrotransposons. PLoS Genet. 8:

706  e1003127.

707 Eichten, S. R., R. A. Swanson-Wagner, J. C. Schnable, A. J. Waters, P. J. Hermanson *et al.*,

708  2011 Heritable epigenetic variation among maize inbreds. PLoS Genet. 7: e1002372.

709 Fu, H., and H. K. Dooner, 2002 Intraspecific violation of genetic colinearity and its implications in

710  maize. Proc. Natl. Acad. Sci. U. S. A. 99: 9573–9578.

711 Gent, J. I., N. A. Ellis, L. Guo, A. E. Harkess, Y. Yao *et al.*, 2013 CHH islands: de novo DNA

712  methylation in near-gene chromatin regulation in maize. Genome Res. 23: 628–637.

713 Haberer, G., N. Kamal, E. Bauer, H. Gundlach, I. Fischer *et al.*, 2020 European maize genomes

714  highlight intraspecies variation in repeat and gene content. Nat. Genet. 52: 950–957.

715 Hirsch, C. N., J. M. Foerster, J. M. Johnson, R. S. Sekhon, G. Muttoni *et al.*, 2014 Insights into

716  the maize pan-genome and pan-transcriptome. Plant Cell 26: 121–135.

717 Hirsch, C. N., C. D. Hirsch, A. B. Brohammer, M. J. Bowman, I. Soifer *et al.*, 2016 Draft

718     Assembly of Elite Inbred Line PH207 Provides Insights into Genomic and Transcriptome

719     Diversity in Maize. Plant Cell 28: 2700–2714.

720  Hoefsloot, H. C., and M. E. Stam, 2020 In plants distal regulatory sequences overlap with

721     unmethylated rather than low-methylated regions, in contrast to mammals. bioRxiv.

722  Hufford, M. B., A. S. Seetharam, and M. R. Woodhouse, 2021 De novo assembly, annotation,

723     and comparative analysis of 26 diverse maize genomes. bioRxiv.

724  Jiao, Y., P. Peluso, J. Shi, T. Liang, M. C. Stitzer *et al.*, 2017 Improved maize reference genome

725     with single-molecule technologies. Nature 546: 524–527.

726  Jiao, Y., P. Peluso, J. Shi, T. Liang, M. C. Stitzer *et al.*, 2016 The complex sequence landscape

727     of maize revealed by single molecule technologies. 1–19.

728  Kawakatsu, T., T. Stuart, M. Valdes, N. Breakfield, R. J. Schmitz *et al.*, 2016 Unique cell-type-

729     specific patterns of DNA methylation in the root meristem. Nature plants 2: 16058.

730  Kim, D., B. Langmead, and S. L. Salzberg, 2015 hisat2. Nat. Methods 944.:

731  Li, H., 2018 Minimap2: pairwise alignment for nucleotide sequences. Bioinformatics 34: 3094–

732     3100.

733  Li, Q., J. I. Gent, G. Zynda, J. Song, I. Makarevitch *et al.*, 2015a RNA-directed DNA methylation

734     enforces boundaries between heterochromatin and euchromatin in the maize genome.

735     Proc. Natl. Acad. Sci. U. S. A. 112: 14728–14733.

736  Li, H., B. Handsaker, A. Wysoker, T. Fennell, J. Ruan *et al.*, 2009 The Sequence

737     Alignment/Map format and SAMtools. Bioinformatics 25: 2078–2079.

738  Li, Q., J. Song, P. T. West, G. Zynda, S. R. Eichten *et al.*, 2015b Examining the causes and

739     consequences of context-specific differential DNA methylation in maize. Plant Physiol. 168:

740     1262–1274.

741  Lister, R., and J. R. Ecker, 2009 Finding the fifth base: Genome-wide sequencing of cytosine

742     methylation. Genome Research 19: 959–966.

743  Love, M., S. Anders, and W. Huber, 2014 Differential analysis of count data--the DESeq2

744      package. Genome Biol. 15: 10–1186.

745      Lu, Z., B. T. Hofmeister, C. Vollmers, R. M. DuBois, and R. J. Schmitz, 2017 Combining ATAC-

746          seq with nuclei sorting for discovery of cis-regulatory regions in plant genomes. Nucleic

747          Acids Res. 45: e41.

748      Marand, A. P., Z. Chen, A. Gallavotti, and R. J. Schmitz, 2020a A cis-regulatory atlas in maize

749          at single-cell resolution.

750      Marand, A. P., Z. Chen, A. Gallavotti, and R. J. Schmitz, 2020b A cis-regulatory atlas in maize

751          at single-cell resolution. bioRxiv.

752      Martin, M., 2011 Cutadapt removes adapter sequences from high-throughput sequencing reads.

753          EMBnet.journal 17: 10–12.

754      Niederhuth, C. E., A. J. Bewick, L. Ji, M. S. Alabady, K. D. Kim *et al.*, 2016 Widespread natural

755          variation of DNA methylation within angiosperms. bioRxiv http://dx.:

756      Noshay, J. M., S. N. Anderson, P. Zhou, L. Ji, W. Ricci *et al.*, 2019 Monitoring the interplay

757          between transposable element families and DNA methylation in maize. PLoS Genet. 15:

758          e1008291.

759      Oka, R., M. Bliek, H. C. J. Hoefsloot, and M. Stam In plants distal regulatory sequences overlap

760          with unmethylated rather than low-methylated regions, in contrast to mammals.

761      Oka, R., J. Zicola, B. Weber, S. N. Anderson, C. Hodgman *et al.*, 2017 Genome-wide mapping

762          of transcriptional enhancer candidates using DNA and chromatin features in maize.

763          Genome Biol. 18: 137.

764      Picard.

765      Regulski, M., Z. Lu, J. Kendall, M. T. Donoghue, J. Reinders *et al.*, 2013 The maize methylome

766          influences mRNA splice sites and reveals widespread paramutation-like switches guided by

767          small RNA. Genome Res. 23: 1651–1662.

768      Ricci, W. A., Z. Lu, L. Ji, A. P. Marand, C. L. Ethridge *et al.*, 2019 Widespread long-range cis-

769          regulatory elements in the maize genome. Nat Plants 5: 1237–1249.

770  Rodgers-Melnick, E., D. L. Vera, H. W. Bass, and E. S. Buckler, 2016 Open chromatin reveals

771     the functional maize genome. Proceedings of the National Academy of Sciences 113:

772     E3177–E3184.

773  Sartor, R. C., J. Noshay, N. M. Springer, and S. P. Briggs, 2019 Identification of the expressome

774     by machine learning on omics data. Proc. Natl. Acad. Sci. U. S. A. 116: 18119–18125.

775  Schmitz, R. J., M. D. Schultz, M. A. Urich, J. R. Nery, M. Pelizzola *et al.*, 2013 Patterns of

776     population epigenomic diversity. Nature 495: 193–198.

777  Schnable, P. S., D. Ware, R. S. Fulton, J. C. Stein, F. Wei *et al.*, 2009 The B73 maize genome:

778     complexity, diversity, and dynamics. Science 326: 1112–1115.

779  Springer, N. M., S. N. Anderson, C. M. Andorf, K. R. Ahern, F. Bai *et al.*, 2018 The maize W22

780     genome provides a foundation for functional genomics and transposon biology. Nat. Genet.

781  Springer, N. M., D. Lisch, and Q. Li, 2016 Creating Order from Chaos: Epigenome Dynamics in

782     Plants with Complex Genomes. Plant Cell 28: 314–325.

783  Springer, N. M., and R. J. Schmitz, 2017 Exploiting induced and natural epigenetic variation for

784     crop improvement. Nat. Rev. Genet. 18: 563–575.

785  Springer, N. M., K. Ying, Y. Fu, T. Ji, C. T. Yeh *et al.*, 2009 Maize inbreds exhibit high levels of

786     copy number variation (CNV) and presence/absence variation (PAV) in genome content.

787     PLoS Genet. 5: e1000734.

788  Stitzer, M. C., S. N. Anderson, N. M. Springer, and J. Ross-Ibarra The Genomic Ecosystem of

789     Transposable Elements in Maize.

790  Sun, S., Y. Zhou, J. Chen, J. Shi, H. Zhao *et al.*, 2018 Extensive intraspecific gene order and

791     gene structural variations between Mo17 and other maize genomes. Nat. Genet. 50: 1289–

792     1295.

793  Swanson-Wagner, R. A., S. R. Eichten, S. Kumari, P. Tiffin, J. C. Stein *et al.*, 2010 Pervasive

794     gene content variation and copy number variation in maize and its undomesticated

795     progenitor. Genome Res. 20: 1689–1699.

796 The Arabidopsis Genome Initiative, 2000 Analysis of the genome sequence of the flowering

797    plant Arabidopsis thaliana. Nature 408: 796–815.

798 Wang, Q., and H. K. Dooner, 2006 Remarkable variation in maize genome structure inferred

799    from haplotype diversity at the bz locus. Proc. Natl. Acad. Sci. U. S. A. 103: 17644–17649.

800 West, P. T., Q. Li, L. Ji, S. R. Eichten, J. Song *et al.*, 2014 Genomic distribution of H3K9me2

801    and DNA methylation in a maize genome. PLoS One 9: 1–10.

802 Xi, Y., and W. Li, 2009 BSMAP: whole genome bisulfite sequence MAPping program. BMC

803    Bioinformatics 10: 232.

804 Xu, J., G. Chen, P. J. Hermanson, Q. Xu, C. Sun *et al.*, 2019 Population-level analysis reveals

805    the widespread occurrence and phenotypic consequence of DNA methylation variation not

806    tagged by genetic variation in maize. Genome Biol. 20: 243.

807 Xu, G., J. Lyu, Q. Li, H. Liu, D. Wang *et al.*, 2020 Evolutionary and functional genomics of DNA

808    methylation in maize domestication and improvement. Nat. Commun. 11: 5539.

809

810

811

812

813

814

815
816
817
818 **Data Tables:**

819
820 **Table 1:** UMR and ACR summary statistics
821

| Sample Genotype | Ref Genotype | # of bins defined as Missing Data | # of bins defined as Methylated* | # of bins defined as Unmethylated | # of UMRs | # of ACRs |
|---|---|---|---|---|---|---|
| B73 | B73v4 | 3511785 | 15064391 | 1325187 | 107178 | 24304 |
| Mo17 | Mo17 | 3649729 | 15566698 | 1385916 | 113838 | 24309 |
| Oh43 | Oh43 | 3096596 | 15719767 | 1445686 | 111261 | 22774 |
| W22 | W22 | 3322802 | 15315985 | 1369207 | 112253 | 21232 |

822 *Methylated is the combined value of bins defined as CG only, CG/CHG, and CHH
823
824
825 **Figure Legends:**
826
827 **Figure 1:** Identification of UMRs and ACRs in maize genotypes.  A) The number of UMRs
828 defined based on samples aligned to B73v4 (green) and their own genome assembly (orange).
829 B) The location of UMRs and ACRs in the genome based on gene annotations was classified as
830 overlapping genes (green), within 2kb of a gene (orange) and >2kb from a gene (purple).   C)
831 The number of ACRs defined based on the merged replicates for each genotype aligned to their
832 respective genome assemblies.  D) Overlap between the B73 UMRs and ACRs defined based
833 on alignments to the B73v4 genome. Number in parentheses indicates ACRs that are defined
834 as methylated as opposed to missing data.
835
836 **Figure 2: Defining shared and nonshared regions between genome assemblies.**  A)
837 Schematic representation of B73-based 100bp bins defined as shared or nonshared in Mo17
838 and W22 (gray shaded regions) based on chromosomal alignments.  The 100bp bins in W22 or
839 Mo17 could be defined by 100bp increments within that genome sequence or based on
840 coordinate matches to the B73 genome and these are shown as the W22 (blue) or Mo17
841 (purple) coordinate bins or the B73-based coordinates (grey).  The black hash or the light to
842 dark color change indicates the 100bp bin boundaries.  B) The proportion of the B73 genome
843 that is defined as shared or non-shared with Mo17, W22, and Oh43 based on chromosome-
844 level sequence alignments.  C) The number of B73 100bp bins that are unique to B73 (0 shared
845 genotypes), shared with one other genotype assessed (1), shared with two other genotypes
846 assessed (2) or shared across all 4 genotypes including B73, Mo17, Oh43, and W22 (3).
847 Genotype labels correspond to the genotypes which share 100bp bins with B73.
848
849 **Figure 3: Presence of ACRs and UMRs within shared and non-shared genomic regions.**
850 A) An IGV (Robinson et al., 2011) representation of a 49kb segment on chromosome 9 of the
851 B73 genome assembly.  Tracks show B73 methylation levels in all contexts (CG-blue, CHG-red,
852 CHH-yellow), B73 UMRs and ACRs, Mo17 shared sequence (green), W22 shared sequence
853 (blue), Oh43 shared sequence (purple), and B73 gene and TE annotations (grey).  B) A small
854 region of the bz1 locus was expanded to see the detail. C) The B73 genome was compared to
855 Mo17, Oh43 or W22 to define regions that are shared or non-shared in each contrast.  The
856 proportion of the shared or non-shared space that is classified as UMR or ACR was determined
857 for each of the pairwise contrasts.
858
859 **Figure 4: Stability of UMRs in shared sequence.**  A) A flowchart on how B73 UMRs are
860 classified is shown.  The numbers in parenthesis indicate the average number of regions
861 classified in that group based on comparisons to the other genotypes. The proportion of B73

862    UMRs that are shared or non-shared (purple) based on sequence with the respective genome
863    assembly.  Shared regions are further classified as B73-only (green) for UMRs that lack data in
864    the other genotype, identical (yellow) for UMRs that maintain an unmethylated state in the same
865    region, partially overlapping (pink) for UMRs that maintain an unmethylated state but have
866    different UMR boundaries across genotypes or polymorphic (blue) for UMRs that change to a
867    methylated state in the other genome. The colors in A are identical to those in C.  B) A genome
868    browser view of the several regions in the B73 genome to illustrate examples of identical,
869    partially overlapping and polymorphic UMRs.  A track of DNA methylation in all contexts (CG-
870    blue, CHG-red, CHH-yellow) is shown for B73 and Mo17 (both aligned to B73v4) with UMRs
871    defined below in black (B73) and blue (Mo17).  B73 UMRs are defined as identical (yellow),
872    partial overlap (pink), or polymorphic (blue).  C) The proportion of B73 UMRs that are classified
873    in each group defined in A are shown for both aUMRs and iUMRs based on comparison to each
874    of the other three genotypes.  D) The number of B73 aUMRs or iUMRs that are classified as
875    ACR only (not unmethylated) in the other genotype (purple), aUMR in the other genotype (blue),
876    iUMR in the other genotype (yellow), or methylated and inaccessible in the other genotype
877    (burgundy) are shown for comparisons to each of the other genotypes

879    **Figure 5: Characteristics of polymorphic UMRs.**  All B73 UMRs classified as polymorphic
880    (shown in Figure 4A) were assessed based on the type of methylation present in the methylated
881    genotype.  The classification is based on which type of methylation state is most common
882    among the 100bp bins of the UMR.  A) A genome browser view of a region on chromosome 5 of
883    the B73 genome.  A track of B73 methylation in all contexts (CG-blue, CHG-red, CHH-yellow) is
884    shown with UMRs defined below in black.  Regions with shared sequence with W22 are shown
885    in red and the W22 methylation track (aligned to the B73v4 assembly) with corresponding UMR
886    classification as overlapping (purple) or polymorphic (red).  Three separate snapshots are
887    shown with the type of methylation found in W22 for the variable UMR noted below (CG only,
888    CG/CHG, or CHH).  B) The percent of all B73 UMRs classified as polymorphic that change to
889    CG only (light blue), CG/CHG (dark blue), or CHH (green) methylation in the other genotype
890    was calculated.  C) UMRs were defined as containing an ACR in both genotypes (Stable ACR:
891    blue), in one genotype (B73 only ACR: green, Non-B73 ACR: orange), or lacking an ACR in
892    both genotypes (No ACR: red).  The proportion of each category of B73 UMR (overlapping and
893    polymorphic) that is defined by ACR presence or absence is shown for each genotype. D) The
894    proportion of UMRs that are found within 200bp of an annotated gene TSS that are defined as
895    differentially expressed (DE), expressed in both genotypes or not expressed is shown for each
896    genotype.  Genes were classified as differentially expressed (log2 fold change > 2 and p-value
897    < 0.05) with the higher expression level observed in B73 (green) or the non-B73 genotype
898    (orange) or as non-differentially expressed (FPKM > 1, pink) or not expressed (silent: purple).

900    **Figure 6: Many differentially methylated tiles (DMTs) are due to partially overlapping**
901    **UMRs.** A) IGV (Robinson et al., 2011) view of DMTs.  Tracks show B73 gene and TE
902    annotations, B73 and Mo17 single cytosine methylation in all contexts (CG: blue, CHG: red,
903    CHH: yellow), B73 UMRs and classification relative to Mo17 (identical: blue, partial: green,
904    polymorphic: red), and DMTs defined by a low level of B73 CG methylation and high level of
905    Mo17 CG methylation.  B) The proportion of B73 DMTs that are associated with partially
906    overlapping UMRs (green) or polymorphic UMRs (orange) is shown. C)  The proportion of B73
907    UMRs, genome-wide (control) or in IBS regions, that are shared or non-shared (purple) based
908    on sequence with the respective genome assembly.  Shared regions are further classified as
909    missing data (orange) for UMRs that lack data in the other genome, identical (blue) for UMRs
910    that maintain an unmethylated state in the same region, partially overlapping (green) for UMRs
911    that maintain an unmethylated state but have different UMR boundaries across genotypes or
912    polymorphic (red) for UMRs that change to a methylated state in the other genome.

913
914
915 **Figure S1:** A) Methylation state classification for each genotype based on alignment to their
916 respective genome assemblies.  Each 100bp bin of the genome was assigned a methylation
917 state based on CG, CHG, and CHH methylation.  Any bin with less than 2 cytosines was labeled
918 "No Sites" and any bin with < 3x coverage was labeled "Missing Data".  For all other bins,
919 context-specific cutoffs of methylation were used to classify CHH, CG only, CG/CHG,
920 Intermediate and Unmethylated status. The proportion of each domain category for all bins in
921 the respective genome are shown.  B) The proportion of all bins for the B73 and non-B73
922 (Mo17, W22, Oh43) samples aligned to the B73v4 genome assembly that have coverage
923 (black) or bins without enough coverage to be assessed (grey).

924
925 **Figure S2**: B73 ATAC-seq reproducibility.  A) Metaplot of ATAC-seq coverage over annotated
926 B73 genes for all ATAC-seq tissue samples aligned to the B73v4 genome assembly.  The gene
927 space was normalized to a 1kb region (represented in the middle of the metaplot) with the
928 flanking upstream and downstream 1kb based on gene transcript direction.  B) ATAC-seq was
929 performed on two replicates for each genotype and ACR calls were generated for each sample
930 individually (BN1A and BN2A) as well as the merged alignment file (BM merge).  The venn
931 diagram represents the overlap in defined ACRs for individual and merged samples for B73.

932
933 **Figure S3:** Overlap between ACRs and UMRs. The overlap between the Mo17 (A), W22 (B)
934 and Oh43 (C) UMRs (blue) and ACRs (green) defined based on alignments to the B73v4
935 genome.  Non-UMR ACRs that are defined as methylated are shown in parentheses below ACR
936 count.  RNAseq data for the same tissue sample was used to classify all B73 genes as not
937 expressed (CPM < 1) or into expression quantiles of lowest expression (Q1) to highest
938 expression (Q4).  For each category of gene, the proportion of genes with aUMRs (D) and
939 iUMRs (E) that are overlapping the gene or proximal to the gene(<2kb) was calculated.

940
941 **Figure S4:** Accessibility is often present only for a portion of the unmethylated region. (A-D)
942 Several B73 UMRs are shown along with ATAC-seq data.  IGV (Robinson et al., 2011)
943 snapshots of the B73 genome showing ACRs within UMR space.  Tracks include B73 gene and
944 TE annotations, B73 methylation per cytosine in all contexts (CG: blue, CHG: red, CHH: yellow),
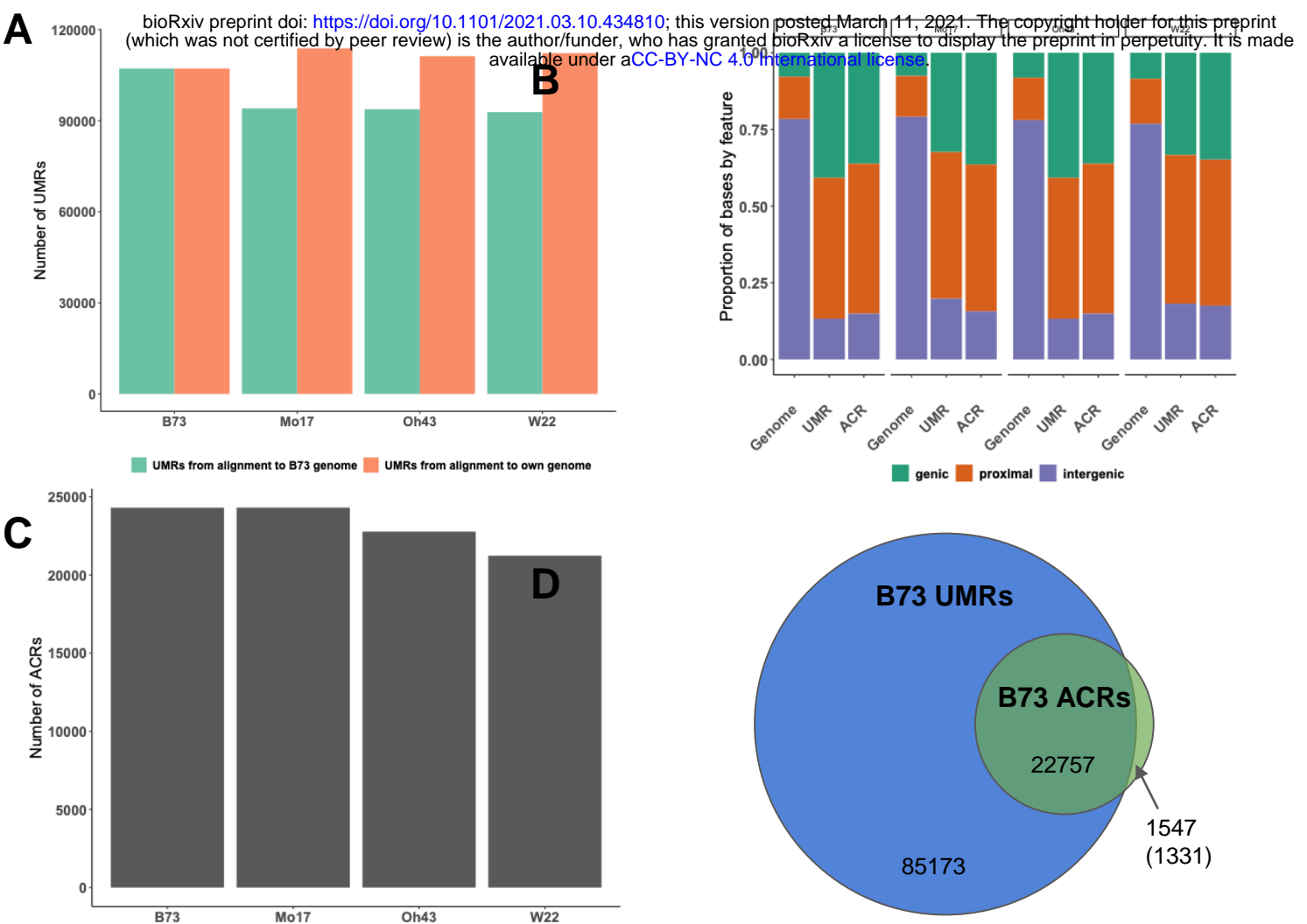945 B73 UMRs (black), B73 ACRs (blue), and B73 ACR coverage (grey).
946

**Figure 1:** Identification of UMRs and ACRs in maize genotypes. A) The number of UMRs defined based on samples aligned to B73v4 (green) and their own genome assembly (orange). B) The location of UMRs and ACRs in the genome based on gene annotations was classified as overlapping genes (green), within 2kb of a gene (orange) and >2kb from a gene (purple). C) The number of ACRs defined based on the merged replicates for each genotype aligned to their respective genome assemblies. D) Overlap between the B73 UMRs and ACRs defined based on alignments to the B73v4 genome. Number in parentheses indicates ACRs that are defined as methylated as opposed to missing data.

**Figure 2: Defining shared and nonshared regions between genome assemblies.** A) Schematic representation of B73-based 100bp bins defined as shared or nonshared in Mo17 and W22 (gray shaded regions) based on chromosomal alignments. The 100bp bins in W22 or Mo17 could be defined by 100bp increments within that genome sequence or based on coordinate matches to the B73 genome and these are shown as the W22 (blue) or Mo17 (purple) coordinate bins or the B73-based coordinates (grey). The black hash or the light to dark color change indicates the 100bp bin boundaries. B) The proportion of the B73 genome that is defined as shared or non-shared with Mo17, W22, and Oh43 based on chromosome-level sequence alignments. C) The number of B73 100bp bins that are unique to B73 (0 shared genotypes), shared with one other genotype assessed (1), shared with two other genotypes assessed (2) or shared across all 4 genotypes including B73, Mo17, Oh43, and W22 (3). Genotype labels correspond to the genotypes which share 100bp bins with B73.
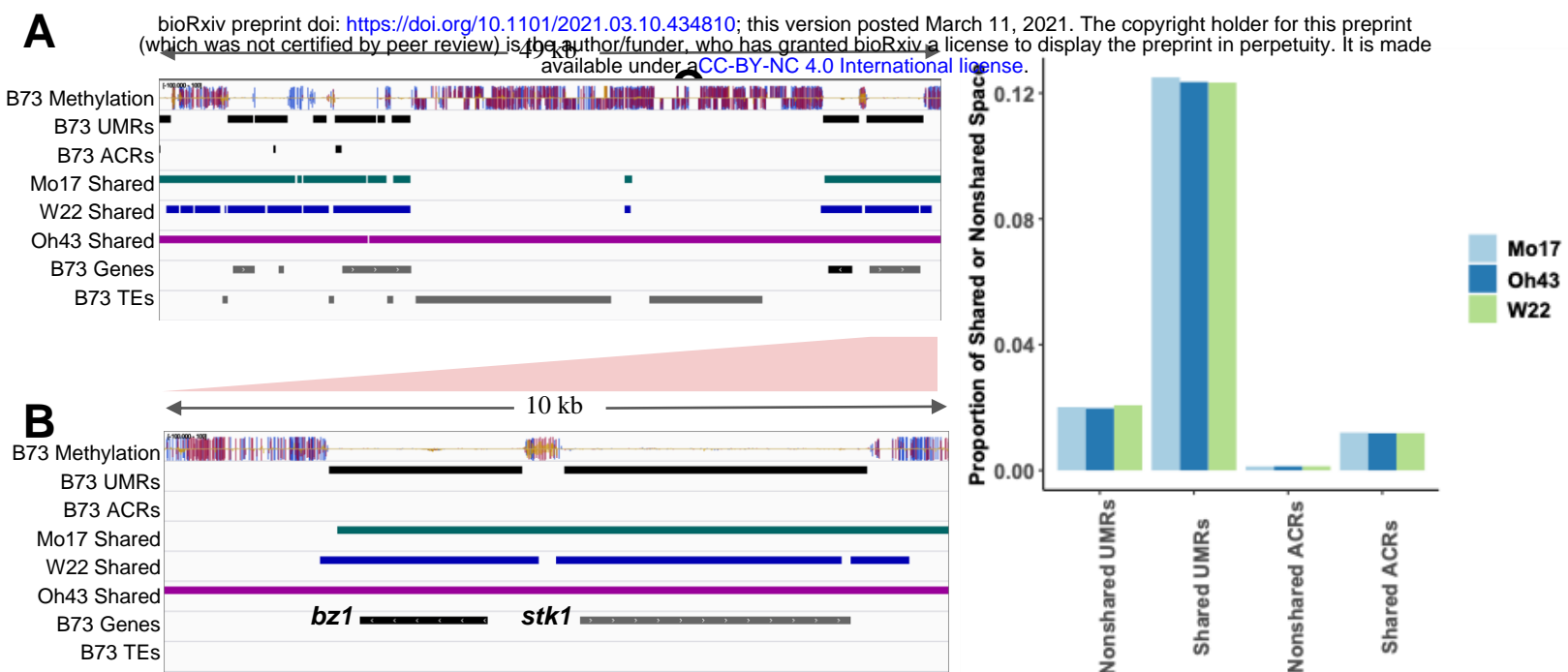
**Figure 3: Presence of ACRs and UMRs within shared and non-shared genomic regions.** A) An IGV (Robinson et al., 2011) representation of a 49kb segment on chromosome 9 of the B73 genome assembly. Tracks show B73 methylation levels in all contexts (CG-blue, CHG-red, CHH-yellow), B73 UMRs and ACRs, Mo17 shared sequence (green), W22 shared sequence (blue), Oh43 shared sequence (purple), and B73 gene and TE annotations (grey). B) A small region of the bz1 locus was expanded to see the detail. C) The B73 genome was compared to Mo17, Oh43 or W22 to define regions that are shared or non-shared in each contrast. The proportion of the shared or non-shared space that is classified as UMR or ACR was determined for each of the pairwise contrasts.

**Figure 4: Stability of UMRs in shared sequence.** A) A flowchart how B73 UMRs are classified is shown. The numbers in parenthesis indicate the average number of regions classified in that group based on comparisons to the other genotypes. The proportion of B73 UMRs that are shared or non-shared (purple) based on sequence with the respective genome assembly. Shared regions are further classified as B73-only (green) for UMRs that lack data in the other genotype, identical (yellow) for UMRs that maintain an unmethylated state in the same region, partially overlapping (pink) for UMRs that maintain an unmethylated state but have different UMR boundaries across genotypes or polymorphic (blue) for UMRs that change to a methylated state in the other genome.The colors in A are identical to those in C. B) A genome browser view of the several regions in the B73 genome to illustrate examples of identical, partially overlapping and polymorphic UMRs. A track of DNA methylation in all contexts (CG-blue, CHG-red, CHH-yellow) is shown for B73 and Mo17 (both aligned to B73v4) with UMRs defined below in black (B73) and blue (Mo17). B73 UMRs are defined as identical (yellow), partial overlap (pink), or polymorphic (blue). C) The proportion of B73 UMRs that are classified in each group defined in A are shown for both aUMRs and iUMRs based on comparison to each of the other three genotypes. D) The number of B73 aUMRs or iUMRs that are classified as ACR only (not unmethylated) in the other genotype (purple), aUMR in the other genotype (blue), iUMR in the other genotype (yellow), or methylated and inaccessible in the other genotype (burgundy) are shown for comparisons to each of the other genotypes
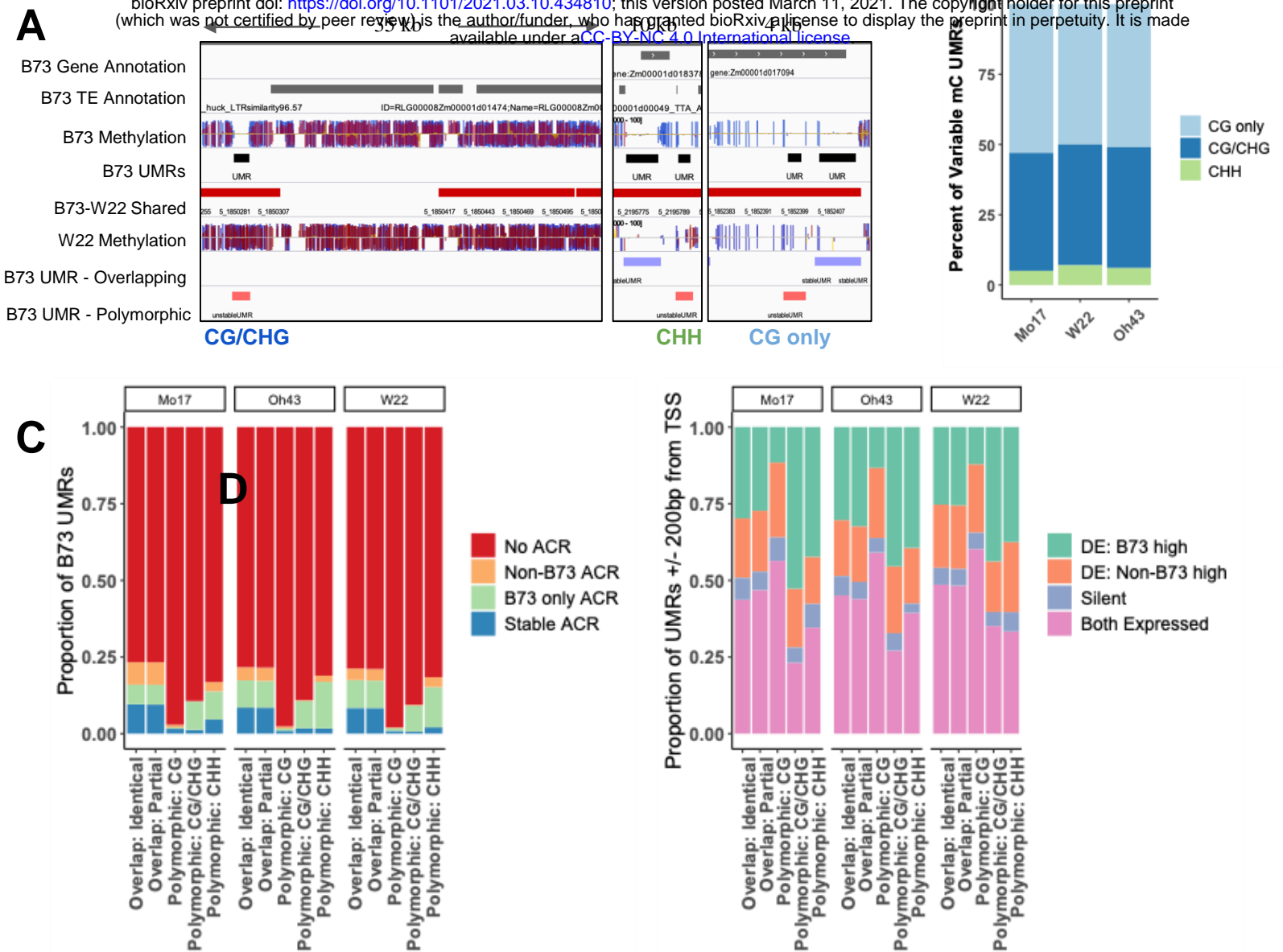
**Figure 5: Characteristics of polymorphic UMRs.** All B73 UMRs classified as polymorphic (shown in Figure 4A) were assessed based on the type of methylation present in the methylated genotype. The classification is based on which type of methylation state is most common among the 100bp bins of the UMR. A) A genome browser view of a region on chromosome 5 of the B73 genome. A track of B73 methylation in all contexts (CG-blue, CHG-red, CHH-yellow) is shown with UMRs defined below in black. Regions with shared sequence with W22 are shown in red and the W22 methylation track (aligned to the B73v4 assembly) with corresponding UMR classification as overlapping (purple) or polymorphic (red). Three separate snapshots are shown with the type of methylation found in W22 for the variable UMR noted below (CG only, CG/CHG, or CHH). B) The percent of all B73 UMRs classified as polymorphic that change to CG only (light blue), CG/CHG (dark blue), or CHH (green) methylation in the other genotype was calculated. C) UMRs were defined as containing an ACR in both genotypes (Stable ACR: blue), in one genotype (B73 only ACR: green, Non-B73 ACR: orange), or lacking an ACR in both genotypes (No ACR: red). The proportion of each category of B73 UMR (overlapping and polymorphic) that is defined by ACR presence or absence is shown for each genotype. D) The proportion of UMRs that are found within 200bp of an annotated gene TSS that are defined as differentially expressed (DE), expressed in both genotypes or not expressed is shown for each genotype. Genes were classified as differentially expressed (log2 fold change > 2 and p-value < 0.05) with the higher expression level observed in B73 (green) or the non-B73 genotype (orange) or as non-differentially expressed (FPKM > 1, pink) or not expressed (silent: purple).
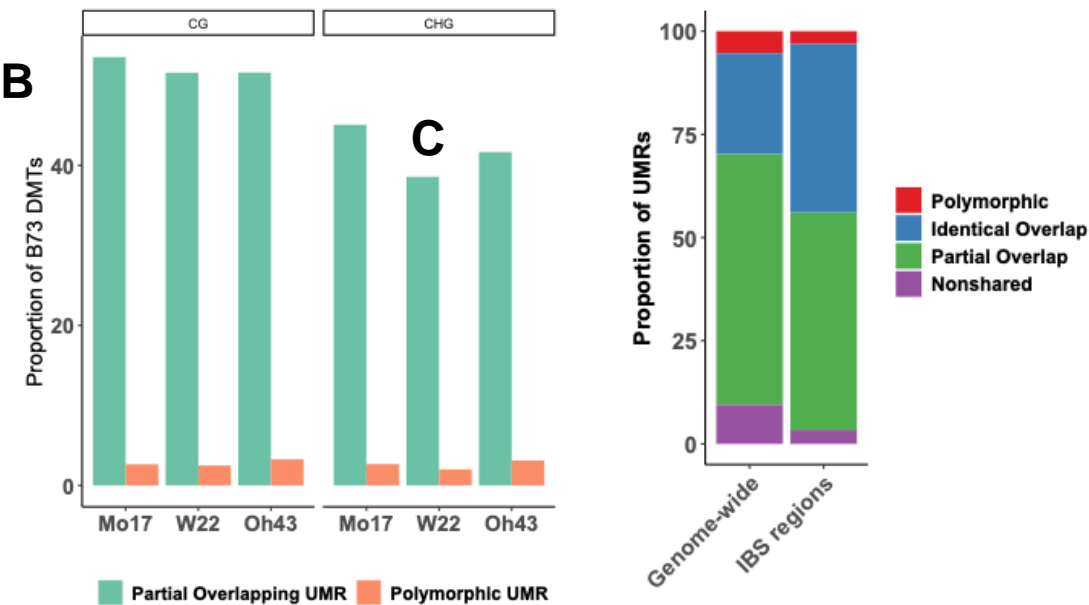
**Figure 6: Many differentially methylated tiles (DMTs) are due to partially overlapping UMRs.** A) IGV (Robinson et al., 2011) view of DMTs.  Tracks show B73 gene and TE annotations, B73 and Mo17 single cytosine methylation in all contexts (CG: blue, CHG: red, CHH: yellow), B73 UMRs and classification relative to Mo17 (identical: blue, partial: green, polymorphic: red), and DMTs defined by a low level of B73 CG methylation and high level of Mo17 CG methylation.  B) The proportion of B73 DMTs that are associated with partially overlapping UMRs (green) or polymorphic UMRs (orange) is shown. C)  The proportion of B73 UMRs, genome-wide (control) or in IBS regions, that are shared or non-shared (purple) based on sequence with the respective genome assembly.  Shared regions are further classified as missing data (orange) for UMRs that lack data in the other genome, identical (blue) for UMRs that maintain an unmethylated state in the same region, partially overlapping (green) for UMRs that maintain an unmethylated state but have different UMR boundaries across genotypes or polymorphic (red) for UMRs that change to a methylated state in the other genome.