

On the Causes of Gene-Body Methylation Variation in *Arabidopsis thaliana*

Rahul Pisupati^{1,2}, Viktoria Nizhynska¹, Almudena Mollá Morales¹, Magnus Nordborg^{1*}

1 Gregor Mendel Institute, Austrian Academy of Sciences, Vienna BioCenter (VBC), Vienna, Austria

2 Vienna Graduate School of Population Genetics, Institut für Populationsgenetik, Vetmeduni, Vienna, Austria

*Corresponding author: magnus.nordborg@gmi.oeaw.ac.at

Abstract

Gene-body methylation (gbM) refers to sparse CG methylation of coding regions, in particular of evolutionarily conserved house-keeping genes. It is found in both plants and animals, but is directly and stably (epigenetically) inherited over multiple generations in the former. Studies in *Arabidopsis thaliana* have demonstrated that plants originating from different parts of the world exhibit large differences in gbM, which presumably reflects an epigenetic memory of ancestral genetic and/or environmental factors.

Here we look for evidence of such factors in F2 plants resulting from a cross between a southern Swedish line with low gbM and a northern Swedish line with high gbM, grown at two different temperatures. Using bisulfite-sequencing data with nucleotide-level resolution on hundreds of individuals, we confirm that CG sites are either methylated (nearly 100% methylation across sampled cells) or unmethylated (approximately 0% methylation across sampled cells), and show that the higher level of gbM in the northern line is due to more sites being methylated. Furthermore, methylation variants almost always show Mendelian segregation, consistent with their being directly and stably inherited through meiosis.

To explore how the differences between the parental lines could have arisen, we focused on somatic deviations from the inherited state, distinguishing between gains (relative to the inherited 0% methylation) and losses (relative to the inherited 100% methylation) at each site in the F2 generation. We demonstrate that deviations predominantly affect sites that differ between the parental lines, consistent with these sites being more mutable. Gains and losses behave very differently in terms of the genomic distribution, and are influenced by the local chromatin state. We find clear evidence for different trans-acting genetic polymorphism affecting both gains and losses, with those affecting gains showing strong environmental interactions (G×E). Direct effects of the environment were minimal.

In conclusion, we show that genetic and environmental factors can change gbM at a cellular level, and hypothesize that these factors can also lead to trans-generational differences between individuals via the inclusion of such changes in the zygote.

Author summary

Gene-body methylation, the sparse CG methylation of house-keeping genes, is found in both plants and animals, but can be directly inherited in the former. Recently, we discovered that *Arabidopsis thaliana* originating from different geographic regions exhibit different patterns of gbM, presumably reflecting a trans-generational memory of genetic or environmental factors. Here we look for evidence of such factors using a genetic cross between two natural inbred lines: one with high, and one with low gbM. We confirm that methylation states are stably inherited, but also see large somatic deviations from the inherited state, in particular at sites that differ between the parental lines. We demonstrate that these deviations are affected by genetic variants in interaction with the environment, and hypothesize that geographic differences in gbM arise through the inclusion of such deviations in the zygote.

Introduction

DNA (cytosine) methylation is an epigenetic mark associated with transcriptional regulation, in particular transposable element silencing. Unlike animals, where methylation is mostly found on CG sites, cytosines in plants are methylated in three-nucleotide contexts: CG, CHG, and CHH, where H=A, C or T. Non-CG methylation is mainly present on transposable elements and is associated with the repression of transcription. It cannot be directly inherited, is found on only a fraction of cells, responds to the environment, and has been shown to be influenced by trans-acting genetic modifiers in *A. thaliana* [1–5]. This is in sharp contrast to CG methylation (mCG), which is maintained during DNA replication through the action of METHYLTRANSFERASE1 (MET1), the homolog of mammalian DNMT1. Unlike in animals, mCHG is not reset every generation in plants, but shows stable trans-generational inheritance [6–9]. As in animals, mCG in plants is present not only on transposable elements and other heterochromatic regions, but also on the coding regions of a subset of genes, a phenomenon known as gene-body methylation (gbM) [10–13]. Genes with gbM tend to be evolutionarily conserved and constitutively expressed, *i.e.*, they are “house-keeping genes”. Although it has been argued that gbM is under selection [14, 15], its function is unclear [16–18].

What is clear is that mCG levels vary greatly between natural inbred lines of *A. thaliana*, and that the pattern of variation reflects the geographic origin of the lines and is correlated with various climate variables [3, 19, 20]. For instance, plants that originate from the colder climate of northern Sweden show higher gbM levels than plants from warmer regions [3]. There are two possible (non-exclusive) explanations for these patterns.

The first is that plants retain an epigenetic memory of their ancestral climate. For this to work, the environment has to affect DNA methylation. Numerous studies have examined the effect of growth conditions on DNA methylation by growing plants in different environments, and while there is clear evidence that non-CG methylation responds strongly to the environment, mCG seems quite stable, at least at the genome-wide level [3, 21–24], consistent with its apparent stability over large numbers of generations [8, 9, 25].

The second is that the geographic pattern of DNA methylation is due to genetic variation. Indeed, genome-wide association studies (GWAS) have identified several trans-acting loci affecting non-CG methylation [3–5, 19, 26], and it is possible that mCG could have been similarly affected by trans-acting modifiers. However, because mCG is stably inherited, it is not a phenotype, and the present methylation state of an individual would not reflect its current genotype but rather the history of its genome, making genetic mapping of such modifiers difficult. It is therefore not surprising that GWAS found no evidence for genetic variants influencing mCG [19].

This paper looks for evidence of genetic variants influencing gbM using a reciprocal F2 cross between a northern Swedish line with high gbM and southern Swedish line with low gbM. To also look for environmental effects, the experiment was carried out at two different temperatures, 4°C and 16°C, and the cross was reciprocal to investigate possible parent-of-origin effects, which are *a priori* plausible [27–29]. Our hope was that our relatively large sample size (a total of over 600 F2 individuals were bisulfite-sequenced) might allow us to detect changes in mCG despite its stable inheritance.

Results

Residual heterozygosity in one parental line

The bisulfite-sequencing data were used to genotype the F2 populations. While doing so, we discovered that one of the parental lines had harbored residual heterozygosity: there are at least two Mb-length regions segregating between the putatively reciprocal F2 populations (S1 Fig). This is irrelevant within each cross, because a single F1 parent was used to generate each F2 population, however, it makes interpretation of differences between the two cross-directions challenging, because they could be due to parent-of-origin effects or genetic differences. For this reason, we will initially focus on the cross in which the northern line was used as mother while the southern was used as father ($n = 308$; S2 Fig), and discuss the (partially) reciprocal cross later. When analyzing parental lines, which were grown in replicate (S1 Table), the segregating regions were eliminated from the analysis.

Differences in gene-body methylation between the parental lines

Methylation estimates from bisulfite sequencing are noisy for a variety of experimental reasons, the most obvious one being low sequence coverage of a possibly heterogeneous population of cells. However, the parental lines were grown in replicate in both temperatures, allowing us to estimate the grandparental state, confirm that gbM is highly consistent between replicates, and that individual sites are either methylated (nearly 100% methylation across sequencing reads) or unmethylated (approximately 0% methylation across sequencing reads), consistent with direct inheritance through both mitosis and meiosis, leading to a cell population with minor deviations from the inherited state, largely independent of temperature (Fig 1, S3 Fig).

The analysis also demonstrated that the previously reported difference in average gbM level between these lines [3] is mostly due to more sites being methylated in the north (rather than a quantitative difference across many sites). Of the roughly 25% of sites that are methylated in at least one of the parental lines, approximately 45% differ between the parental lines, and, of these, 70% are only methylated in the northern line (Fig 1).

Inheritance of gene-body methylation in the F2 population

In the F2 population we do not have replication of entire genotypes, but we have massive replication of the genotype at each site, because 1/4 of the 308 F2 individuals are expected to be homozygous for northern ancestry (NN) at each site, 1/4 to be homozygous for southern ancestry (SS), and 1/2 to be heterozygous (NS). Ancestry can accurately be inferred using SNP haplotypes, and by combining this with the methylation states in the F2 population we can also infer the epigenotype at each site in the F1 parent — and confirm that gbM shows the expected Mendelian segregation (S3 Fig, [30]).

The inferred F1 epigenotype can be compared with the inferred grand-parental epigenotype to get an estimate of the epimutation rate. This is not straightforward and requires a number of assumptions because differences could have arisen at any point across two generations — and could also reflect heterozygosity in the grand-parental generation, as well as various artefacts that are difficult to control for. We obtain a per-generation, per-site rate of loss of gbM of $\sim 0.2\%$, and corresponding rate of gain of $\sim 0.04\%$, although we caution that there are aspects of our data we cannot explain (see Materials and Methods for details).

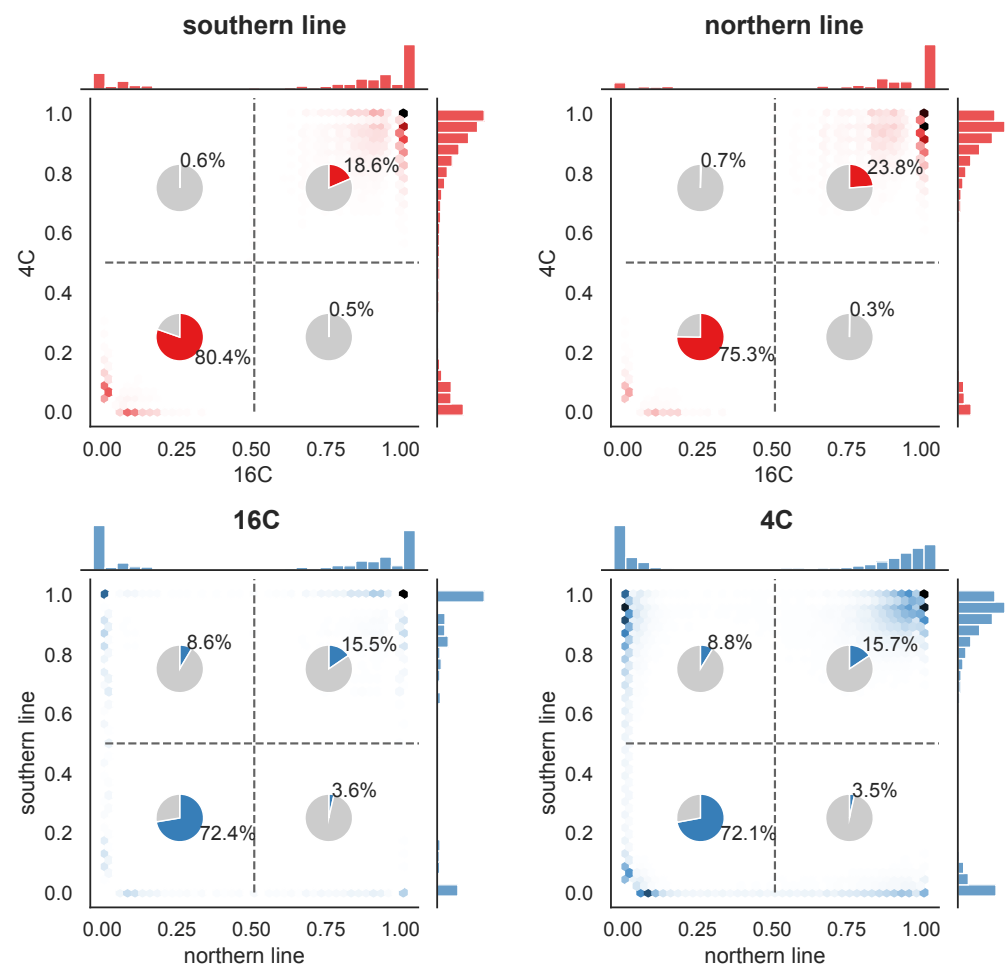


Fig 1. The pattern of gbM across sites. The plots show the distribution of average methylation levels across 650,595 gbM sites at 4°C and 16°C, separately for the two parental lines. The pie charts show the fraction of sites classified as methylated or unmethylated using 50% methylation as a cut-off (see Materials and Methods for details). The top plots compare temperatures for each parental line; the bottom plots compare parental lines for each temperature.

However, these epimutations did not occur in the F2 generation. While they may have been affected by the F1 genotype, they do not reflect genetic variants segregating in the F2 population, nor our temperature treatment. In order to take advantage of the experimental design, we need to focus on changes that happened in the F2 generation itself, *i.e.*, we need a proper phenotype. Thus we focus on somatic deviations from the inherited state (as seen in the parents in Fig 1). These are by definition phenotypes affected by genotype and environment, and while the deviation at a particular site in a particular individual is very poorly estimated (primarily due to insufficient sequencing coverage), this is compensated by the size of the F2 population. It is obvious from Fig 1 that *gains* (positive deviations from an inherited state of 0% methylation) have a very different distribution from *losses* (negative deviations from an inherited state of 100% methylation), and we therefore estimate each separately, as explained in Fig 2.

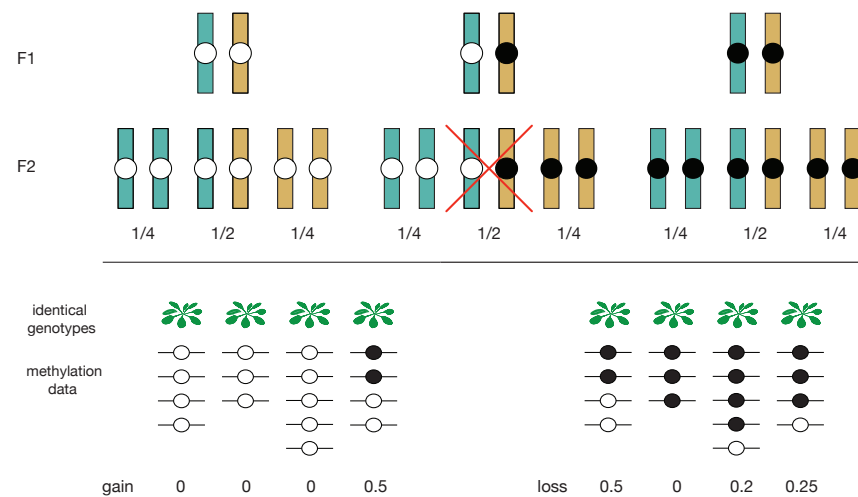


Fig 2. Quantifying somatic gains and losses. In the F2 population, each gbM site is characterized by ancestry: NN, NS, and SS. Independently of this, there are three types of sites: those for which F1 parent was homozygous unmethylated, those for which it was heterozygous methylated/unmethylated (could be either on N or S allele), and those for which it was homozygous methylated. In the F2 population we estimate gains only for individuals that should have inherited the homozygous unmethylated state, and losses only for individuals that should have inherited the homozygous methylated state. We do not use individuals heterozygous for methylation. Different analyses then use different subsets of the gain/loss data as detailed below.

Losses and gains reflect different processes

Somatic losses and gains differ in several aspects. First, estimated losses are on average two orders of magnitude higher than estimated gains: 7.3% vs 0.09%, respectively (S2 Table). Second, gains and losses show very different distributions across the genome, similar to what has been observed for trans-generational epimutations [31]. Gains are 2.8 times higher in peri-centromeric regions, while losses are correlated with the enrichment of the H2A.Z chromatin mark on the gene ($r = 0.14$; $p < 0.01$) (Fig 3). Third, losses vary much more across the four possible CG contexts (CGA, CGT, CGC, CGG) than gains. In particular, losses are 22% higher on CGT compared to the other contexts (S4 Fig).

Gains and losses are only weakly affected by temperature (Fig 3, S4 Fig). They do, however, depend on local ancestry: on average losses are 2% higher on SS alleles compared to NN alleles, while gains are higher 29% on NN alleles than SS alleles (although the pattern varies greatly across the genome; see Fig 3). Potential causes for these patterns will be discussed below. Finally, both gains and losses exhibited positive auto-correlation along the genome (gains are correlated with gains at nearby sites, and the same for losses). We do not observe any such effects on non-CG methylation (S4 Fig).

Importantly, both gains and losses are higher for sites that differ between the two parental lines: the increase is roughly 10-fold for gains and almost 2-fold for losses (Fig 3D, S3 Table). Given this, and the other similarities to trans-generational epimutations noted above, it is reasonable to hypothesize that both are generated by the same mechanism, and that trans-generational epimutations are simply somatic epimutations that end up being transmitted via gametes.

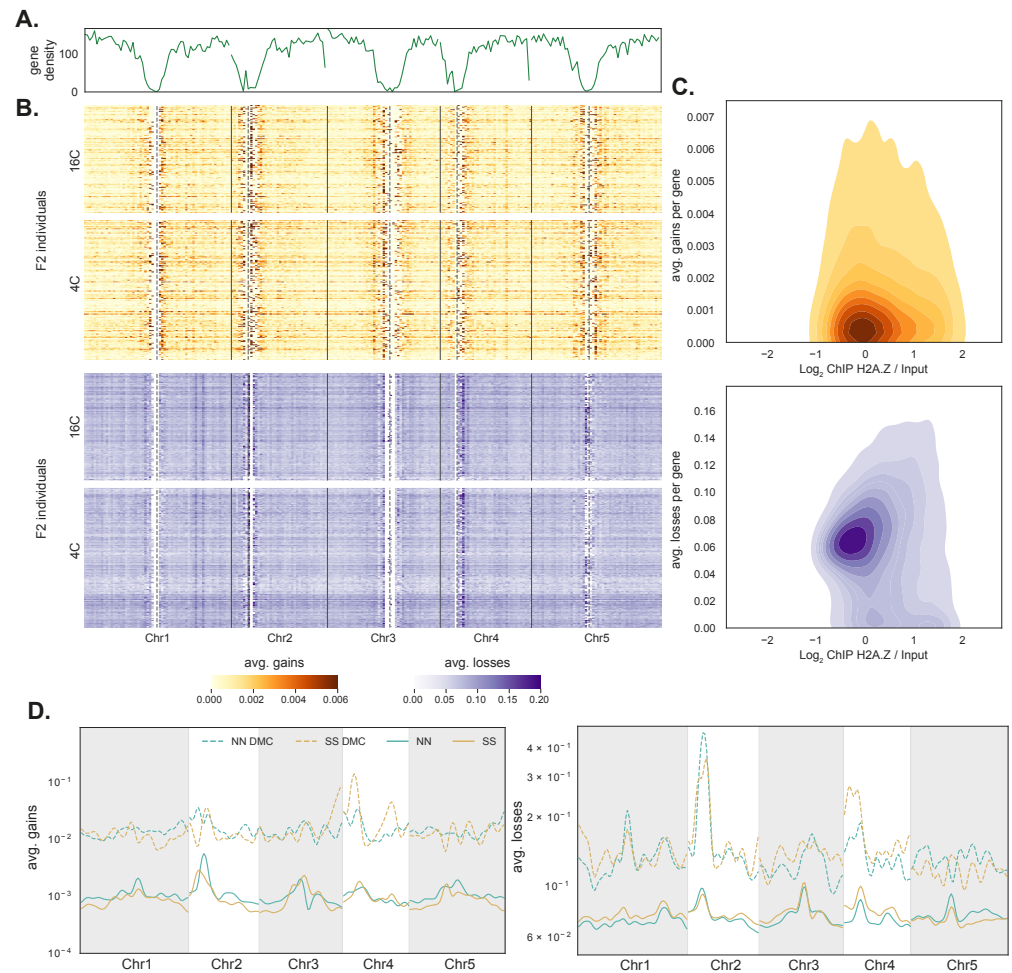


Fig 3. Somatic deviations across the genome (A) Line plot for gene density. (B) Heatmaps show genome-wide somatic deviations for gains and losses for genes in F2 individuals at both temperatures (n=308). Each row is an individual. Gene density and deviations were calculated in 500 kb windows across the genome. Vertical solid lines represent chromosome breaks and dotted lines represent the centromere positions. (C) Average gains and losses for each gene (in NN background) plotted against average H2A.Z ChIP-seq level (in Col-0 [32]). (D) Average gains and losses across the genome for homozygous NN and SS individuals. Deviations at sites that differ between the parents are shown using dashed lines (see Fig 2).

Motivated by this, we investigated whether the observed gains and losses have the properties one would naïvely expect of mitotically heritable epimutations. Specifically, we tested whether cells switch state independently of each other (conditional on estimated rates of switching) within and between individuals using a simulation approach (see Materials and Methods). If deviations were largely due to somatically inherited epimutations (perhaps effecting large sectors of the sequenced plants), changes within plants would be positively correlated, and we might see inflated variance between plants, with some plants being responsible for most of the average deviation at a given site (see Fig 2). However, with the possible exception of gains on sites that differ between the parents, we see no evidence of this phenomenon (S5 Fig). The distribution of gains seems compatible with independent changes within and

between plants, and there is no evidence for large sectors due to somatic inheritance (*n.b.* our power to detect such sectors is extremely limited due to low sequencing coverage per-individual).

The distribution of losses, on the other hand, is clearly incompatible with independent mutations, but in the opposite direction: there is far too little variation between individuals for losses to reflect random independent events (S5 Fig).

Genetic architecture of deviations

To investigate genetic and environmental factors influencing these deviations, we used a standard F2 linkage mapping model that includes temperature as an environmental factor and allows for genotype-by-environment interaction ($G \times E$). As phenotypes, we used deviations in 500 kb windows across the genome. Windows were used because per-site deviations are far too noisy (since deviations are rare), and using genome-wide deviations is inappropriate given clear evidence for heterogeneity across the genome (Fig 3): the 500 kb size was empirically determined. The results provide further evidence that gains and losses are different phenomena. For both phenotypes, we identify significant *trans*-acting QTL, but they are not the same (Fig 4A, S4 Table). Furthermore, gains are also affected by strong *cis*-acting factors.

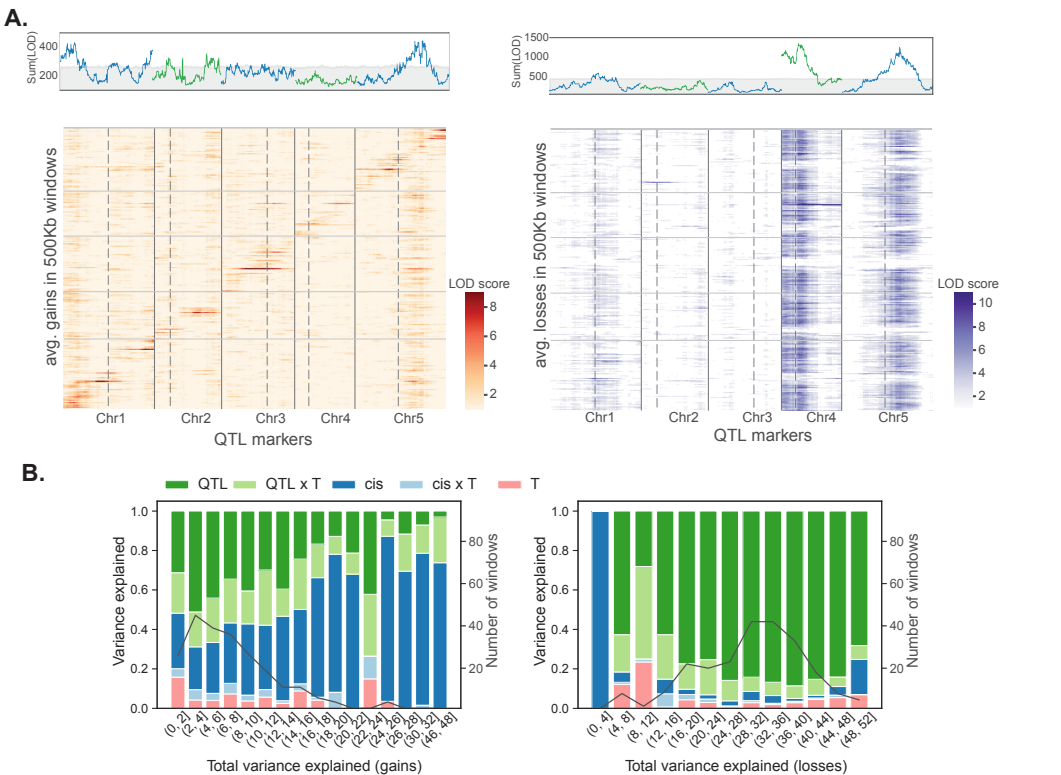


Fig 4. Genetic architecture of deviations. (A) Heatmaps showing linkage mapping results for gains and losses in 500 kb windows together with plots summing LOD scores across these windows. Peaks above gray region are significant using a 1% FDR based on genome permutations. Vertical dotted lines identify centromeres and solid lines separate chromosomes. (B) Bar plots summarizing variance partitioning results for gains and losses. Results are binned by total variance explained, with thin black lines showing the distribution of windows across bins.

QTL affecting losses are far stronger and had more consistent effects across the genome. We identify two major QTL accounting for about 5% of the variation each, with similar effects in both temperatures, and with additive effects within and between loci (*i.e.*, no dominance or epistasis; see Fig 5, S6 Fig and S4 Table). The northern and southern alleles have opposite effects at the two loci.

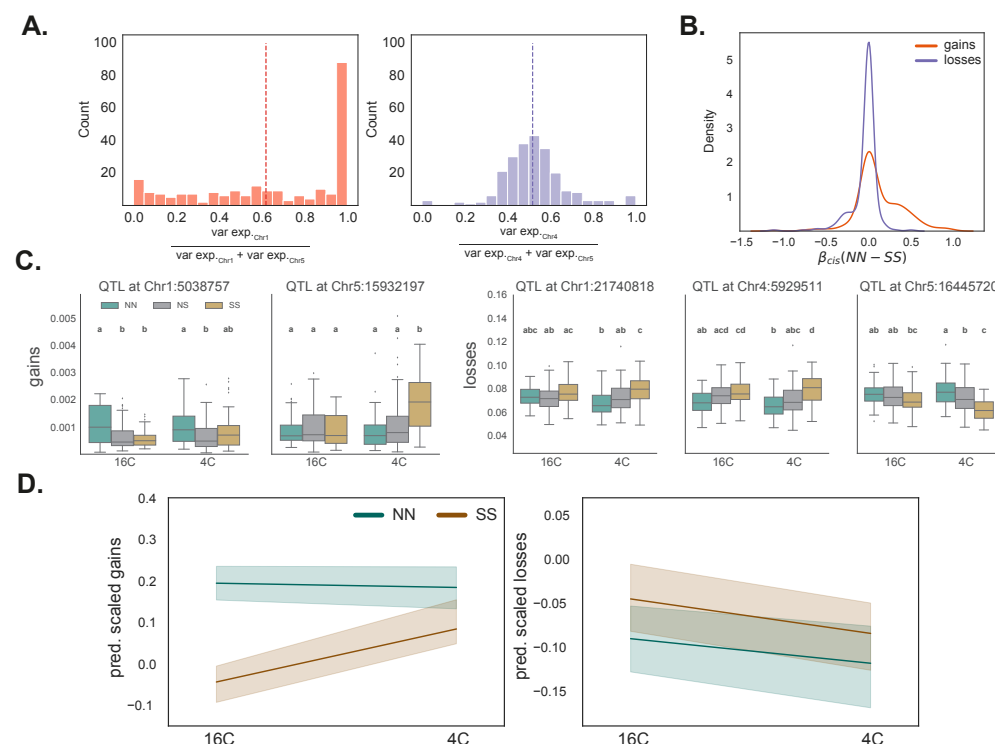


Fig 5. QTL effect-size estimates. (A) The distribution of variance explained across 500 kb windows for a gain QTL (left) and a loss QTL (right). The mean effects (vertical lines) are similar, but the gain QTL has a highly skewed distribution, with strong effect only on a subset of windows, where the loss QTL affects most of the genome. (B) The distribution of *cis* effect sizes. (C) Genotypic effects for two gain QTL and three loss QTL. Average gains and losses across windows significantly associated with QTL are shown. (D) Reaction norms for predicted gains and losses in individuals homozygous for the northern or southern alleles at all significant QTLs.

The two QTL for gains affect different windows (Fig 5). Each QTL explains a couple of percent of the variation, and the north-south direction of effects is again reversed between the loci. At each locus, the allele associated with greater gains is recessive, and the effect of the chromosome 5 QTL is only seen at 4°C. There is no evidence for epistasis.

In order to quantify the factors affecting the deviations, we partitioned the variance in each 500 kb window using a linear model that includes local (*cis*-) genotype (*i.e.*, NN, SS, or NS), temperature, and the identified QTL as factors. The results for gains and losses are again strikingly different (Fig 4B). For losses, the majority of the variance is explained by the QTL, with a minor role for QTL-by-temperature (QTL×T) interactions. For gains, QTL, QTL×T, and *cis*-genotype appear to play roughly equal roles, and there is also evidence for interactions between the *cis*-genotype and temperature. Temperature, in-and-of-itself,

explains little of the variation, however, the G×E effects for gains are substantial. This can also be seen in the predicted response for the parental QTL genotypes (Fig 5D), which agree with direct estimates (S7 Fig).

In an attempt to fine-map some of the QTL identified here, we turned to GWAS. We used the population data from reference [3], where about 100 accessions were grown at two temperatures, 10°C and 16°C. We calculated genome-wide deviations for each accessions by considering sites with less than 50% methylation as gains and sites with more than 50% methylation as losses. Consistent with temperature having little effect, deviations are highly correlated between the two temperatures ($r = 0.74$, S9 Fig). The average gains and losses across accessions are around 0.5% and 9%, respectively, and we performed GWAS using these as phenotypes, but could not identify any significant associations (S9 Fig). The same is true when 500 kb windows rather than genome-wide averages are used.

Cis-effects on deviations

We have seen that deviations are associated with the local haplotype, *i.e.*, they are affected by *cis*-acting factors (Fig 4). The effect is particularly pronounced for gains, but is also seen for losses. Generally speaking, the *cis*-effects work in the direction of the observed differences, *i.e.*, gains are more pronounced on the more methylated N allele and losses are higher on the less methylated S allele (Fig 5B).

While it possible that these effects are due to genetics, it would imply that *cis*-regulatory differences have evolved throughout the genome. It seems more likely that the effects are a consequence of the epigenetic differences that we know exist. As mentioned previously, deviations are associated with the underlying chromatin state (Fig 3), suggesting the local epigenetic state influence them.

Motivated by this, we examined whether deviations are correlated with methylation levels at the level of individual genes. And indeed, gains tended to be higher on the allele with higher methylation level, while losses show the opposite pattern (Fig 6A). Zooming in further, we find that both gains and losses are strongly affected by nearby methylation at a nucleotide scale (Fig 6B). For gains in particular, the effect seems to be limited to less than 30 bp.

Partially reciprocal cross

As noted above, this experiment was designed to include a reciprocal cross in order to test for parent-of-origin effects on methylation, but undetected residual heterozygosity in one of the parental lines made the cross only partially reciprocal, making interpretation of differences challenging. For this reason, discussion thus far has been limited to the cross in which the northern line was used as maternal parent.

In the reciprocal cross, we observe very similar patterns of deviations across the genome (S8 Fig). Average losses are strongly correlated between the F2 populations at the level of genes (Fig 7A), and the two significant QTL appear to be replicated (although the significance of the one on chromosome 5 was weaker). In addition, we identify a new QTL on chromosome 1 that directly overlapped the region segregating between the F2 populations and is thus probably due to a genetic difference rather than the direction of the cross (Fig 7B).

In stark contrast, average gains are not correlated between the two directions of the cross (Fig 7A). Given this, it is not surprising that the corresponding QTL mapping results are also discordant, with previously identified QTL being replaced by different ones in (Fig 7B). The QTL do not overlap regions that segregate between the F2 populations, and must thus either reflect epistatic interaction with putative causal polymorphism in these regions, or parent-of-origin effects.

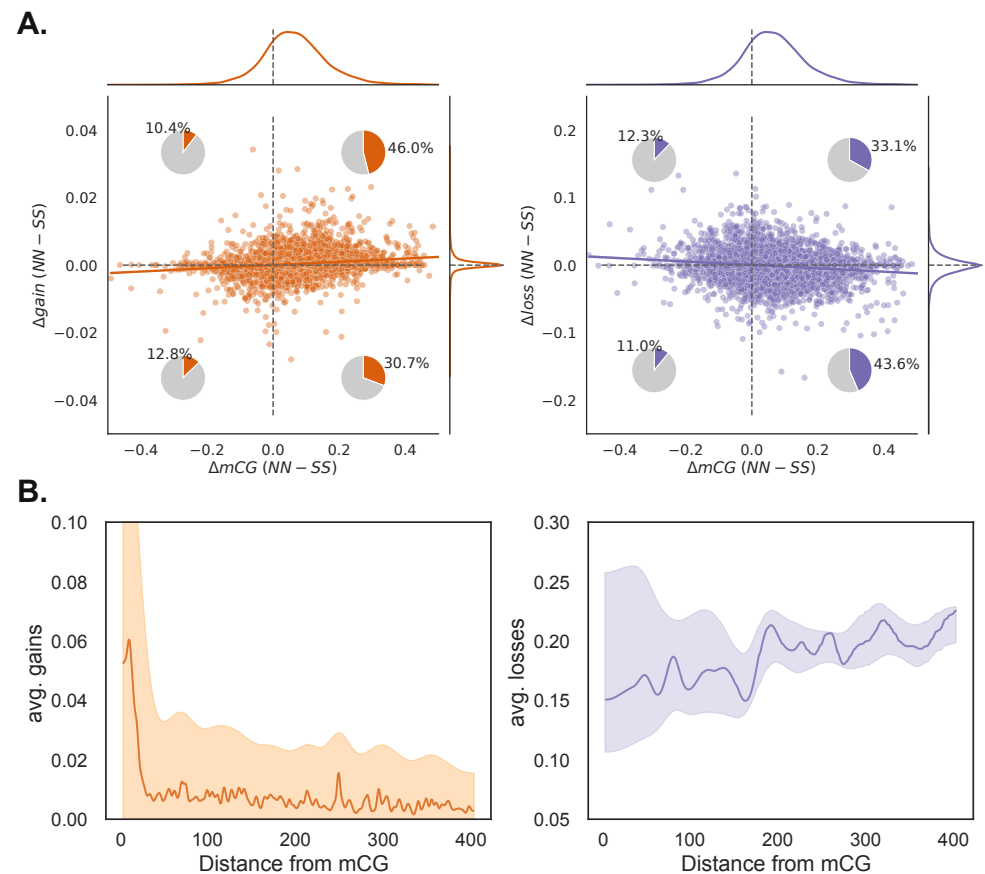


Fig 6. *Cis*-effects on deviations (A) Across 10,160 genes with gbM, the difference in gains between northern and southern alleles (estimated from homozygous individuals) is correlated with the difference in gbM between the same alleles. Both correlations (Spearman coefficients of 0.2 and -0.12 for gains and losses respectively) are significant ($p < 0.01$). **(B)** Average gains and losses at a given CG site depends on the distance to the nearest methylated CG site.

Discussion

The motivation for this study was to gain insight into how gbM is inherited — and how it changes. While several studies have established that mCG is generally stably inherited over large numbers of generations, albeit with a high (epi-)mutation rate [1, 8, 9, 33], it is also clear that substantial geographic differences exist, differences that cannot be explained via random mutations [3, 19]. We used a traditional diallel F2 cross between two parental lines that differ considerably in gbM to investigate this further. Our analysis provides very strong confirmation that mCG shows Mendelian segregation [30], and our estimated per-site, per-generation epimutation rates of $\sim 0.04\%$ for gain of methylation and $\sim 0.2\%$ for loss of methylation are also consistent with previous estimates (although there are odd phenomena that will be discussed below) [33].

What is novel about our study is that we focus on somatic deviations from the inherited methylation state, either gains (for sites inherited as unmethylated) or losses (for sites inherited as methylated). Not only do these provide more observations than trans-generational epimutations (since we survey more cells than plants), but, more

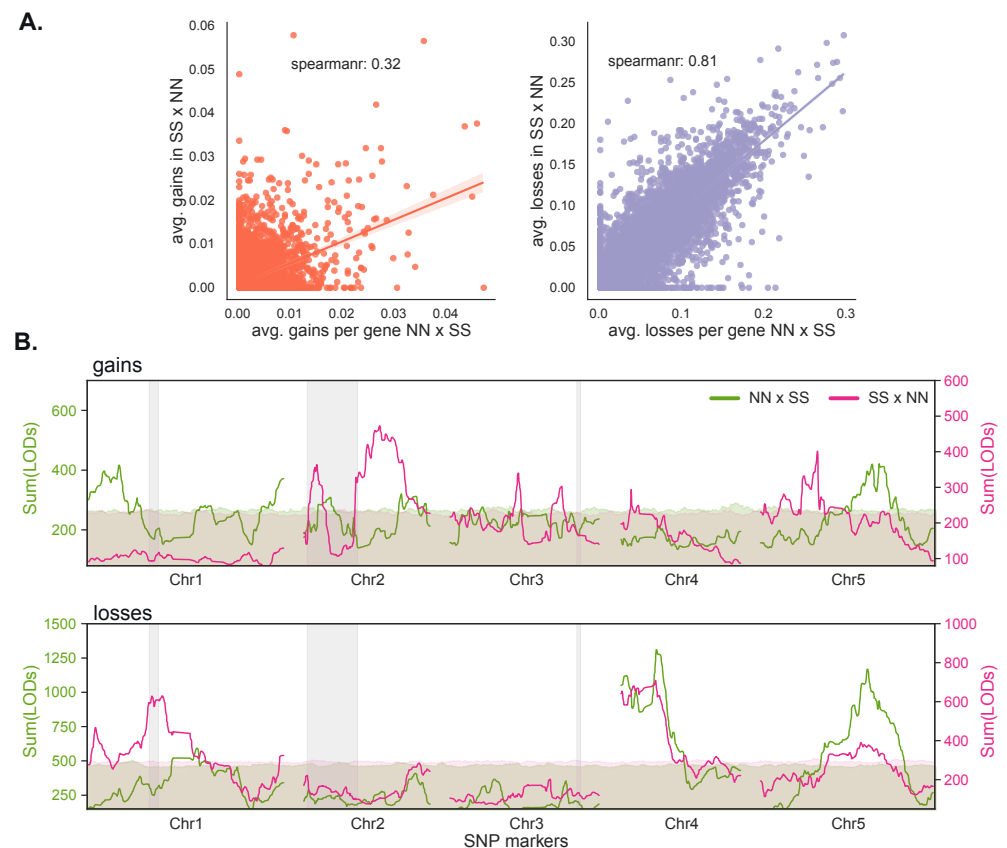


Fig 7. Somatic deviations in the partially reciprocal cross. (A) Average gains and losses per gene (in NN background) for both directions of the cross. (B) Aggregate linkage mapping results (sum of LOD scores for 500 kb window) for both directions of the cross. The results for the NN×SS direction were already shown in Fig 4A. The 99% significance thresholds were determined using 1000 genome rotations (see Materials and Methods). Regions that segregate between the two F2 populations are shown using grey vertical bands.

importantly, they are by definition phenotypes — they occurred in the current generation, and could have been affected by the genetic background and environmental exposure of each individual. These effects can be investigated using standard methods of quantitative genetics.

We find that gains and losses behave very differently, presumably reflecting different molecular mechanisms. Gains occur at low rates (higher than the estimated trans-generational rate of gains, but same order of magnitude), independently within and between individuals, perfectly consistent with their being somatic epimutations. We also see that somatic gains of methylation are positively correlated with nearby (within 30 bp) methylation, just as has been observed for trans-generational gains of methylation (*cf.* Fig 6B with Fig 3 in reference [34]). It is thus reasonable to hypothesize that the QTL we identify correspond to *bona fide* modifiers of the epimutation rate — which makes it very interesting that they show strong G×E effects (Fig 4), as well as possible parent-of-origin effects (Fig 7). If the mechanisms that give rise to the somatic gains we observe also give rise to trans-generational gains, the pattern of gbM variation observed in nature [3,19] could reflect a complex interplay between trans-acting genetic factors and the environment.

Somatic losses, on the other hand, at first look nothing like epimutations. They occur at rates two orders-of-magnitude higher than trans-generational epimutations, and are furthermore positively correlated between individuals, which is clearly not consistent with random mutations. One explanation is that they reflect experimental artefacts due to bisulfite sequencing, which typically shows less than 100% methylation when used on fully methylated control DNA [35]. However, while this is likely to contribute, artefacts would not give rise to highly significant QTL. These must have a biological basis, but not necessarily one related to epimutations. It is important to remember that mCG is automatically lost during DNA replication (the newly synthesized strand is unmethylated, leading to hemi-methylated DNA), and that the maintenance of mCG across mitosis is therefore an active process, catalyzed by MET1 [1]. Anything that caused an imbalance between the rate of cell division and MET1 activity could lead to somatic losses, and these could well be unrelated to trans-generational epimutations.

This said, the somatic losses we observe are probably not completely unrelated to epimutations. At least two lines of evidence speak against this. First, both losses and gains are much more pronounced on sites that differ between the parental lines (Fig 3), consistent with a shared mutational mechanism. Second, just as was the case for gains, losses show a dependence on local methylation that is similar to what has been seen for trans-generational loss-of-methylation mutations (*cf.* Fig 6B with Fig 3 in reference [34]).

Considering all this, we believe that the losses we observe reflect a mixture of (biased) experimental noise and biological factors that are distinct from those affecting gains. In addition to the differences in the fine-scale pattern also observed in reference [34], the QTL we identify for losses have larger effect than those for gains, and show no evidence of G×E or parent-of-origin effects.

In conclusion, we have shown that while gbM methylation is apparently mostly directly inherited, it can be influenced by trans-acting genetic modifiers that are different for gains and losses, and that show strong G×E effects for the former. Whether such modifiers can explain the natural geographic variation in gbM remains to be demonstrated, but is a plausible hypothesis. Finally, we emphasize that much remains unknown about gbM. We do not understand its function (if any), and we do not even fully understand how it is established and maintained. For the latter question, better data (*e.g.*, strand-specific methylation haplotype data from single cells not obtained using bisulfite-sequencing) in large pedigrees will be needed.

Materials and methods

Plant growth

We chose two natural inbred lines from Sweden that had been shown to differ considerable in gbM [3]: one line from Lövvik in northern Sweden (ID 6046, lat. 62.800323, long. 18.075722) with average gbM of 16% and another from Drakamöller in southern Sweden (ID 6191, lat. 55.758856, long. 14.132822) with average gbM of 12.5%. We generated recombinant hybrid progeny of these two lines by collecting seeds from selfed F1 individuals for the reciprocal directions (S1 Fig). Selfed parental lines were grown along with F2 individuals from two families at two temperatures (16°C and 4°C) in a randomized block design (S1 Table). We grew plants on soil and stratified for 5 days at 4°C in the dark before transferring them to long day chambers with 16 hours of light and 8 hours of darkness. When plants attained the 9-true-leaf stage of development, one or two leaves were collected and flash-frozen in liquid nitrogen.

DNA extraction and bisulfite sequencing

Genomic DNA was extracted from frozen tissue using the NuclearMag Plant kit (Machery-Nagel). We adopted a tagmentation-based protocol to generate multiplexed DNA libraries for whole-genome bisulfite sequencing (T-WGBS; [35]). We optimized the protocol for low DNA inputs (20 ng) and high-throughput (96-well plates). We used in-house Tn5 transposase generated at Vienna BioCenter Core Facilities. The tagmentation, oligonucleotide replacement and gap repair were done according to the T-WGBS protocol.

We used EZ-96 DNA Methylation-Gold Mag Prep kit (Zymo Research) for bisulfite conversion from tagmented DNA. We PCR-amplified bisulfite-treated DNA with 15-16 cycles with KAPA HiFi Uracil polymerase (Kapa Biosystems). We used Illumina TruSeq unique index adapters for PCR amplification and multiplexing of the libraries. Amplified libraries were validated using Fragment Analyzer™ Automated CE System (Advanced Analytical) and multiplexed (96X) in equimolar concentration. Libraries were sequenced on Illumina HiSeq™2000 Analyzers or HiSeqV4 using the manufacturer's standard cluster generation and sequencing protocols in 100-125bp paired-end mode.

Sequencing data analysis

Sequenced BS-seq reads were analyzed using a well-documented nf-core pipeline (github.com/nf-core/methylseq). First, BS-seq reads were trimmed for adaptors using cutadapt (default parameters), and we clipped 15 bp at the beginning of the reads due to uneven base composition. Second, the trimmed reads were mapped to the TAIR10 (Col-0) reference assembly using bismark relaxing mismatches to 0.5 [36]. Third, methylation calling was performed using methylpy on the aligned bam files. We used custom scripts to calculate weighted averages of methylation [37] at annotated genes and transposable elements using the ARAPORT11 annotation (www.arabidopsis.org/download/index-auto.jsp?dir=%2Fdownload_files%2FGenes%2FAraport11_genome_release). All scripts used were packaged in python and are available on github (github.com/Gregor-Mendel-Institute/pyBsHap.git).

Bisulfite conversion rate estimation

It is common practice to use the chloroplast genome (cpDNA) to estimate conversion rates for BS-seq libraries, since cpDNA is unmethylated [38]. The non-conversion rate was calculated as the fraction of methylated cytosines from reads mapped to the cpDNA. The estimated conversion rate for the libraries is on average 0.3%. We then ignored methylation on sites that did not have significantly higher methylation than expected due to non-conversion (using a binomial test with probability of 0.3%; p-value of 0.05).

Gene body methylation

We calculated methylation levels on all exonic CG sites. We excluded genes with significant non-CG methylation in either of the parental lines from the analysis (S10 Fig), but did not rely on any other epigenetic marks. In doing so, the average mCHG and mCHH levels per gene were scaled, and outlier genes were identified using twice the standard deviation.

SNP calling and genetic map reconstruction using bisulfite treated libraries

The mapped bam files from bismark were modified for base positions that are influenced by bisulfite treatment ($C \rightarrow T$ and $G \rightarrow A$) using Revelio ([github.com/bio15anu/revelio.git](https://github.com/bio15anu/revelio)) [39]. We genotyped 10.7 million previously identified SNP sites [40] using bcftools with default parameters [41]. The scripts for the pipeline were packaged and hosted on github (github.com/Gregor-Mendel-Institute/nf-haplocaller).

Next, we inferred the underlying ancestry at each SNP marker segregating between parents in F2 individuals using a multinomial hidden Markov model (adapted from reference [42]) packaged in the SNPmatch package (github.com/Gregor-Mendel-Institute/SNPmatch.git). Bisulfite sequencing gives uneven coverage across the genome, but such data can be used to infer ancestry with high accuracy, in particular in F2 individuals where ancestry tracts are very long. We filtered out SNP markers having identical genotype data across individuals using R/qtl package [43]. This resulted in a total of 3983 SNP markers used for linkage mapping.

Residual heterozygosity in reciprocal cross

We calculated percentage of heterozygous SNP calls for parental lines sequenced as part of the 1001 Genomes project [40]. At least four genomic regions more than 300 kb had residual heterozygosity in the southern parent (S1 FigB).

As a consequence, for any given site in these regions, different southern alleles could be segregating in the reciprocal crosses, *i.e.*, rather than N and S alleles segregating in both, we could have N and S_1 in one direction, and N and S_2 in the other. To identify such regions, we identified all SNP segregating in each F2 population, then compared them using SNPmatch [44]. As expected, this revealed that a subset of the putatively heterozygous regions differed between the directions of the cross (S1 Fig).

Estimating somatic deviations

Each F2 family ($NN \times SS$ and $SS \times NN$) is the offspring of a single F1 individual, a hybrid with NS-ancestry at every site. Every mCG site would either be methylated (11), unmethylated (00) or heterozygous (01) in this F1 individual (Fig 2). Due to the stable inheritance of mCG, we expect the parental methylation state to have been passed on, and this was readily confirmed. Somatic gains and losses were calculated as weighted averages across sites classified as having been inherited homozygous unmethylated or methylated, respectively [37]. This was either done per gene or in windows of 500 kb.

In individuals heterozygous for methylation state (NS), we expect to see 50% methylation since we lack the power to do allele-specific methylation (given 100 bp reads, and our data supports this (S3 Fig).

The python scripts used for these analyses were packaged and are hosted on github (github.com/Gregor-Mendel-Institute/pyBsHap).

Modeling somatic deviations

Let s_{ij} be the number of reads with ancestral methylation at site i in individual j , and let f_{ij} be the number of reads with non-ancestral methylation. We calculate deviation from the ancestral state as $x_{ij} = f_{ij}/n_{ij}$, where $n_{ij} = s_{ij} + f_{ij}$. We also define the

average deviation at site i ,

$$\bar{x}_{i.} = \sum_{j=1}^N x_{ij}/N;$$

the average deviation for individual j ,

$$\bar{x}_{.j} = \sum_{i=1}^M x_{ij}/M;$$

and the total average deviation

$$\bar{x} = \sum_{i=1}^M \bar{x}_{i.}/M = \sum_{j=1}^N \bar{x}_{.j}/N.$$

We wish to test the null-model that deviations are due to independent mutation in each cell, mutations that occur with site- and individual-specific probabilities. For site i in individual j , reads were simulated by drawing from a binomial distribution with parameters n_{ij} and $p_{ij} = \bar{x}_{i.} + \bar{x}_{.j}$. We then calculated the variance across individuals for each site, and compared simulation results with data. If there were large sectors of epimutations in some individuals (*i.e.*, non-independence of states within individuals), the between-individual variance should be inflated. We observe the opposite for losses, whereas gains are broadly consistent with the null model.

QTL mapping and variance partitioning

We performed linkage mapping using the R/qtl package [43]. We use both simple interval mapping (using the ‘scanone’ function) and composite interval mapping (using the ‘cim’ function) via Haley-Knott regression. We included growth temperature as a cofactor when performing linkage mapping as full model. We identified QTLs having an interaction with temperature by comparing full model with the additive model. QTLs were identified by adding LOD scores across genomic regions. The significance threshold was calculated by permuting ($n = 1000$) LOD scores and performing genome-rotations to retain the LD structure.

We estimated variance explained for identified QTLs, *cis* genotype, temperature, and their interactions using a linear mixed model. We used the ‘lmer’ function from ‘lme4’ package in R [45] to implement the model

$$y = G_{cis} + T + \sum_i G_{QTL_i} + G_{cis} \times T + \sum_i G_{QTL_i} \times T + \epsilon, \quad (1)$$

where y is the somatic deviation at a given genomic region, G_{cis} is the genotype at the *cis* marker, G_{QTL_i} is the genotype at QTL marker i , and T is the growth temperature.

Genome-wide association mapping (GWAS)

GWAS was performed using a linear mixed model implemented using LIMIX [46]. We used the SNP matrix ($n = 3,916,814$) from the 1001 Genomes Project filtered for SNPs with minor allele frequency greater than 5% in the Swedish populations [40].

Estimating epimutation rates

We used the average methylation across replicate individuals of each parental line to infer the methylation state of the grand-parental individual (Fig 1 and S1 Fig).

Analogously, we used a weighted average across F2 individuals homozygous for each ancestry (NN or SS) to infer the methylation state of the F1 individual (separately for the N and S chromosomes, see S1 Fig).

Any trans-generational epimutations that occurred during these two generations would give rise to differences between the inferred grand-parental and F1 states, and it should be possible to use this to estimate the epimutation rate. However, such differences could also result from estimation error, and we realized that one important source of such error would be sites heterozygous for methylation in the grand-parental generation. Such sites would lead to segregating methylation among the parental replicates, and lead to random assignment of grand-parental methylation state using our 50% rule. They are expected to be extremely rare, and indeed there is no evidence of them in Fig 1. However, when comparing the distribution of methylation levels across sites in the averaged parental individuals compared to the averaged F2 individuals, we do see an enrichment of sites with intermediate methylation in the former (S11 Fig), presumably reflecting grand-parental heterozygosity. Another potential source of error is cryptic copy number variation, which could lead to pseudo-heterozygosity and again intermediate levels of methylation [47].

In order to guard against these errors, we filtered out all sites with ambiguous methylation state in either the grand-parental or F1 generation, conservatively retaining only sites consistent with the genome-wide distributions of gains and losses (average somatic gains per site are less than 0.2 and average somatic losses per site are less than 0.35, see S5 Fig). Almost all sites removed using this approach are due to insufficient coverage in the parental generation (which does not affect the calculations of somatic gains and losses in the rest of the paper).

With this filtered data, we estimate epimutation rates by comparing the inferred methylation state of the grandparent with that of the F1. Any difference effectively means that the F1 allele must have changed state either in early development, or via an epimutation from parent to F1, or from grandparent to parent, *i.e.*, the changes reflect two generations.

We calculated epimutation rates separately for the northern and southern lineage, and also for the two F1 individuals resulting from the two directions of the reciprocal cross (S12 Fig, S13 Fig, S5 Table). The average per-site, per-generation epimutation rates are $\sim 0.04\%$ and $\sim 0.2\%$ respectively (S5 Table), but there are several anomalies that caution against over-interpretation of these estimates. First, losses on the northern lineage are three times higher in the NN \times SS direction of the cross than in the SS \times NN direction, and gains on the southern lineages are two times higher in the NN \times SS direction of the cross than in the SS \times NN direction. Second, when filtering for ambiguous methylation the F1 generation, we detected evidence of rapid change in this generation, consistent with the action of trans-acting modifiers. Third, the overlap in mutated sites between the two crosses is orders of magnitude higher than could be expected under any model of random mutations. The far greater sharing along the northern lineage suggests that the same parental individual was used as mother in one direction of the cross and father in the other (this would result in sharing of half of all epimutations that occurred in the first of the two generations of the pedigree, see S13 Fig). These observations provide further evidence (see also reference [34]) that a model of random epimutations is not sufficient, and suggest that further experiments are badly needed.

Finally, we calculated epimutation rates for sites differentially methylated between the parental lines (S14 Fig, S6 Table). Consistent with the patterns in somatic deviations (Fig 3), epimutation rates are much higher for these sites: ten-fold for gains and two-fold for losses.

Data availability

The raw sequencing data and the methylation calls are uploaded to NCBI GEO database under GSE215839. All the scripts and the intermediate data files used for the analysis are uploaded on Github (github.com/Gregor-Mendel-Institute/pisupati-gbm-paper-2022.git).

Acknowledgments

We are grateful to Ortrun Mittelsten Scheid, Fred Berger, and Kelly Swarts for discussions throughout the project. We would also like to thank the Nordborg lab, especially Yoav Voichkek, Haijun Liu and Thomas Ellis for their helpful comments on the manuscript. Finally, we thank Bob Schmitz and Daniel Zilberberman for comments on this paper and Rahul Pisupati's thesis. Bisulfite sequencing was performed by the Next Generation Sequencing Facility at the Vienna BioCenter Core Facilities (VBCF), a member of the Vienna BioCenter (VBC), Austria. This research was supported by ERC AdG 789037 "EPICLINES" to MN.

References

1. Law JA, Jacobsen SE. Establishing, maintaining and modifying DNA methylation patterns in plants and animals. *Nature Reviews Genetics*. 2010;11(3):204–220.
2. Kawashima T, Berger F. Epigenetic reprogramming in plant sexual reproduction. *Nature Reviews Genetics*. 2014;15(9):613–624.
3. Dubin MJ, Zhang P, Meng D, Remigereau MS, Osborne EJ, Casale FP, et al. DNA methylation in Arabidopsis has a genetic basis and shows evidence of local adaptation. *elife*. 2015;4:e05255.
4. Sasaki E, Kawakatsu T, Ecker JR, Nordborg M. Common alleles of CMT2 and NRPE1 are major determinants of CHH methylation variation in Arabidopsis thaliana. *PLoS genetics*. 2019;15(12):e1008492.
5. Sasaki E, Gunis J, Reichardt-Gomez I, Nizhynska V, Nordborg M. Conditional GWAS of non-CG transposon methylation in Arabidopsis thaliana reveals major polymorphisms in five genes. *PLoS genetics*. 2022;18(9):e1010345. doi:10.1371/journal.pgen.1010345.
6. Saze H, Mittelsten Scheid O, Paszkowski J. Maintenance of CpG methylation is essential for epigenetic inheritance during plant gametogenesis. *Nature genetics*. 2003;34(1):65–69. doi:10.1038/ng1138.
7. Reinders J, Wulff BB, Mirouze M, Marí-Ordóñez A, Dapp M, Rozhon W, et al. Compromised stability of DNA methylation and transposon immobilization in mosaic Arabidopsis epigenomes. *Genes & development*. 2009;23(8):939–950.
8. Becker C, Hagmann J, Müller J, Koenig D, Stegle O, Borgwardt K, et al. Spontaneous epigenetic variation in the Arabidopsis thaliana methylome. *Nature*. 2011;480(7376):245–249.
9. Schmitz RJ, Schultz MD, Lewsey MG, O'Malley RC, Urich MA, Libiger O, et al. Transgenerational epigenetic instability is a source of novel methylation variants. *Science*. 2011;334(6054):369–373.

10. Tran RK, Henikoff JG, Zilberman D, Ditt RF, Jacobsen SE, Henikoff S. DNA methylation profiling identifies CG methylation clusters in Arabidopsis genes. *Current biology: CB*. 2005;15(2):154–159. doi:10.1016/j.cub.2005.01.008. 514
515
516
11. Zilberman D, Gehring M, Tran RK, Ballinger T, Henikoff S. Genome-wide analysis of Arabidopsis thaliana DNA methylation uncovers an interdependence between methylation and transcription. *Nature genetics*. 2007;39(1):61–69. doi:10.1038/ng1929. 517
518
519
520
12. Lister R, O'Malley RC, Tonti-Filippini J, Gregory BD, Berry CC, Millar AH, et al. Highly integrated single-base resolution maps of the epigenome in *Arabidopsis*. *Cell*. 2008;133(3):523–536. 521
522
523
13. Bewick AJ, Schmitz RJ. Gene body DNA methylation in plants. *Current opinion in plant biology*. 2017;36:103–110. 524
525
14. Takuno S, Gaut BS. Body-methylated genes in Arabidopsis thaliana are functionally important and evolve slowly. *Molecular biology and evolution*. 2012;29(1):219–227. doi:10.1093/molbev/msr188. 526
527
528
15. Muyle AM, Seymour DK, Lv Y, Huettel B, Gaut BS. Gene Body Methylation in Plants: Mechanisms, Functions, and Important Implications for Understanding Evolutionary Processes. *Genome biology and evolution*. 2022;14(4). doi:10.1093/gbe/evac038. 529
530
531
532
16. Zilberman D. An evolutionary case for functional gene body methylation in plants and animals. *Genome biology*. 2017;18(1):1–3. 533
534
17. Bewick AJ, Ji L, Niederhuth CE, Willing EM, Hofmeister BT, Shi X, et al. On the origin and evolutionary consequences of gene body DNA methylation. *Proceedings of the National Academy of Sciences*. 2016;113(32):9111–9116. 535
536
537
18. Sarda S, Zeng J, Hunt BG, Yi SV. The evolution of invertebrate gene body methylation. *Molecular biology and evolution*. 2012;29(8):1907–1916. 538
539
19. Kawakatsu T, Huang SSC, Jupe F, Sasaki E, Schmitz RJ, Urlich MA, et al. Epigenomic diversity in a global collection of Arabidopsis thaliana accessions. *Cell*. 2016;166(2):492–505. 540
541
542
20. Schmitz RJ, Schultz MD, Urlich MA, Nery JR, Pelizzola M, Libiger O, et al. Patterns of population epigenomic diversity. *Nature*. 2013;495(7440):193–198. 543
544
21. Vaughn MW, Tanurdzić M, Lippman Z, Jiang H, Carrasquillo R, Rabinowicz PD, et al. Epigenetic natural variation in Arabidopsis thaliana. *PLoS biology*. 2007;5(7):e174. 545
546
547
22. Downen RH, Pelizzola M, Schmitz RJ, Lister R, Downen JM, Nery JR, et al. Widespread dynamic DNA methylation in response to biotic stress. *Proceedings of the National Academy of Sciences*. 2012;109(32):E2183–E2191. 548
549
550
23. Eichten SR, Springer NM. Minimal evidence for consistent changes in maize DNA methylation patterns following environmental stress. *Frontiers in plant science*. 2015;6:308. 551
552
553
24. Wibowo A, Becker C, Marconi G, Durr J, Price J, Hagmann J, et al. Hyperosmotic stress memory in Arabidopsis is mediated by distinct epigenetically labile sites in the genome and is restricted in the male germline by DNA glycosylase activity. *Elife*. 2016;5:e13546. 554
555
556
557

25. Hagmann J, Becker C, Müller J, Stegle O, Meyer RC, Wang G, et al. Century-scale Methylome Stability in a Recently Diverged *Arabidopsis thaliana* Lineage. *PLoS genetics*. 2015;11(1):e1004920. doi:10.1371/journal.pgen.1004920.
26. Hüther P, Hagmann J, Nunn A, Kakoulidou I, Pisupati R, Langenberger D, et al. MethylScore, a pipeline for accurate and context-aware identification of differentially methylated regions from population-scale plant whole-genome bisulfite sequencing data. *Quantitative Plant Biology*. 2022;3:e19. doi:10.1017/qpb.2022.14.
27. Hsieh TF, Shin J, Uzawa R, Silva P, Cohen S, Bauer MJ, et al. Regulation of imprinted gene expression in *Arabidopsis* endosperm. *Proceedings of the National Academy of Sciences*. 2011;108(5):1755–1762.
28. Gehring M. Genomic imprinting: insights from plants. *Annual review of genetics*. 2013;47:187–208. doi:10.1146/annurev-genet-110711-155527.
29. Satyaki PRV, Gehring M. DNA methylation and imprinting in plants: machinery and mechanisms. *Critical reviews in biochemistry and molecular biology*. 2017;52(2):163–175. doi:10.1080/10409238.2017.1279119.
30. Hofmeister BT, Lee K, Rohr NA, Hall DW, Schmitz RJ. Stable inheritance of DNA methylation allows creation of epigenotype maps and the study of epiallele inheritance patterns in the absence of genetic variation. *Genome biology*. 2017;18(1):1–16.
31. Hazarika RR, Serra M, Zhang Z, Zhang Y, Schmitz RJ, Johannes F. Molecular properties of epimutation hotspots. *Nature plants*. 2022;doi:10.1038/s41477-021-01086-7.
32. Jamge B, Lorković Z, Axelsson E, Yelagandula R, Akimcheva S, Berger F. Transcriptional activity is shaped by the chromatin landscapes in *Arabidopsis*. *bioRxiv*. 2022;.
33. Van Der Graaf A, Wardenaar R, Neumann DA, Taudt A, Shaw RG, Jansen RC, et al. Rate, spectrum, and evolutionary dynamics of spontaneous epimutations. *Proceedings of the National Academy of Sciences*. 2015;112(21):6676–6681.
34. Briffa A, Hollwey E, Shahzad Z, Moore JD, Lyons DB, Howard M, et al. Unified establishment and epigenetic inheritance of DNA methylation through cooperative MET1 activity. *bioRxiv*. 2022;doi:10.1101/2022.09.12.507517.
35. Wang Q, Gu L, Adey A, Radlwimmer B, Wang W, Hovestadt V, et al. Tagmentation-based whole-genome bisulfite sequencing. *Nature protocols*. 2013;8(10):2022–2032.
36. Krueger F, Andrews SR. Bismark: a flexible aligner and methylation caller for Bisulfite-Seq applications. *Bioinformatics*. 2011;27(11):1571–1572. doi:10.1093/bioinformatics/btr167.
37. Schultz MD, Schmitz RJ, Ecker JR. 'Leveling' the playing field for analyses of single-base resolution DNA methylomes. *Trends in genetics: TIG*. 2012;28(12):583–585. doi:10.1016/j.tig.2012.10.012.
38. Schmitz RJ, Marand AP, Zhang X, Mosher RA, Turck F, Chen X, et al. Quality control and evaluation of plant epigenomics data. *The Plant cell*. 2022;34(1):503–513. doi:10.1093/plcell/koab255.

39. Nunn A, Otto C, Fasold M, Stadler PF, Langenberger D. Manipulating base
quality scores enables variant calling from bisulfite sequencing alignments using
conventional Bayesian approaches. *BMC genomics*. 2022;23(1):1–10. 602
603
604
40. The 1001 Genomes Consortium. 1,135 Genomes Reveal the Global Pattern of
Polymorphism in *Arabidopsis thaliana*. *Cell*. 2016;166(2):481–491. 605
doi:10.1016/j.cell.2016.05.063. 606
607
41. Li H. A statistical framework for SNP calling, mutation discovery, association
mapping and population genetical parameter estimation from sequencing data.
Bioinformatics. 2011;27(21):2987–2993. doi:10.1093/bioinformatics/btr509. 608
609
610
42. Andolfatto P, Davison D, Erezyilmaz D, Hu TT, Mast J, Sunayama-Morita T,
et al. Multiplexed shotgun genotyping for rapid and efficient genetic mapping.
Genome research. 2011;21(4):610–617. doi:10.1101/gr.115402.110. 611
612
613
43. Broman KW, Wu H, Sen S, Churchill GA. R/qtl: QTL mapping in experimental
crosses. *Bioinformatics*. 2003;19(7):889–890. doi:10.1093/bioinformatics/btg112. 614
615
44. Pisupati R, Reichardt I, Seren Ü, Korte P, Nizhynska V, Kerdaffrec E, et al.
Verification of *Arabidopsis* stock collections using SNPmatch, a tool for
genotyping high-plexed samples. *Scientific data*. 2017;4(1):1–9. 616
617
618
45. Bates D, Mächler M, Bolker B, Walker S. Fitting Linear Mixed-Effects Models
Using lme4. *Journal of Statistical Software*. 2015;67(1):1–48. 619
doi:10.18637/jss.v067.i01. 620
621
46. Lippert C, Casale FP, Rakitsch B, Stegle O. LIMIX: genetic analysis of multiple
traits. *BioRxiv*. 2014; p. 003905. 622
623
47. Jaegle B, Pisupati R, Soto-Jiménez LM, Burns R, Rabanal FA, Nordborg M.
Extensive gene duplication in *Arabidopsis* revealed by pseudo-heterozygosity.
bioRxiv. 2022;doi:10.1101/2021.11.15.468652. 624
625
626

Supporting information

627

S1 Table. Individuals sequenced.

Genotype	16°C	4°C	Total
Parent N	7	16	23
Parent S	9	16	25
F2 (NN×SS)	134	174	308
F2 (SS×NN)	132	174	306

S2 Table. Average somatic deviations in F2 individuals.

	Gains	Losses
16°C	0.00092	0.0744
4°C	0.001	0.0724
Mean	0.00097	0.073

S3 Table. Deviations in NN and SS backgrounds, separately for sites that are identical vs differ between N and S.

	Gains		Losses	
	Identical	Differ	Identical	Differ
NN	0.0011	0.0152	0.0738	0.1428
SS	0.00085	0.01470	0.0753	0.1558

S4 Table. Linkage mapping results for deviations in NN×SS cross.

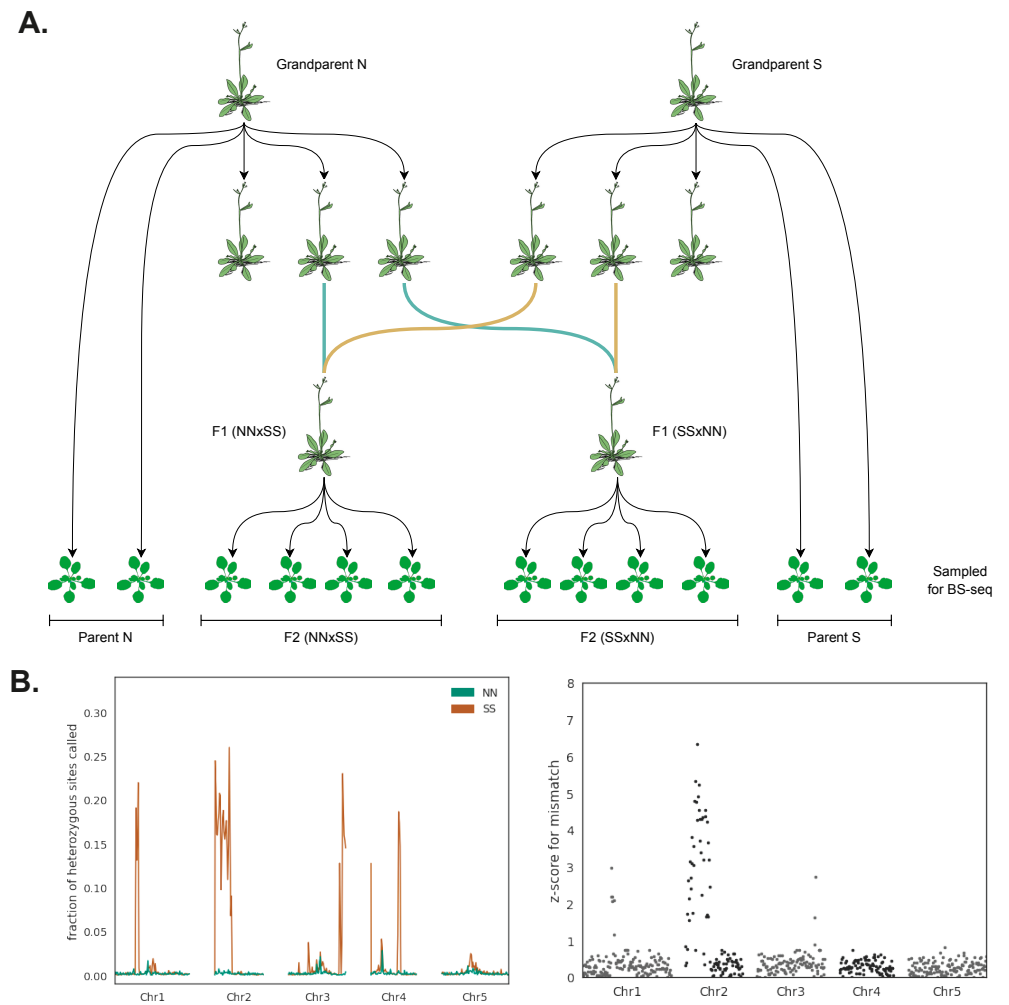
	Top SNP	95% CI in Mb	Candidate genes (position)
Gains	Chr1:5038757	Chr1:4.0-5.5	RDR1 (5.1), SHH1 (5.2), IDNL1 (5.4), SUVH7 (6.1)
	Chr5:15932197	Chr5:15.8-17.6	VIM3 (15.8), AGO10 (17.6), SUVR2 (17.7)
Losses	Chr1:21740818	Chr1:8.0-23.6	VIM1 (21.4), NRPD1 (23.3)
	Chr4:5929511	Chr4:5.1-6.3	RDR2 (6.8), SUVH9 (7.8), MET2 (8.1)
	Chr5:16445720	Chr5:16.1-17.0	VIM3 (15.8), AGO10 (17.6), SUVR2 (17.7)

S5 Table. Epimutation rates using data from S12 Fig.

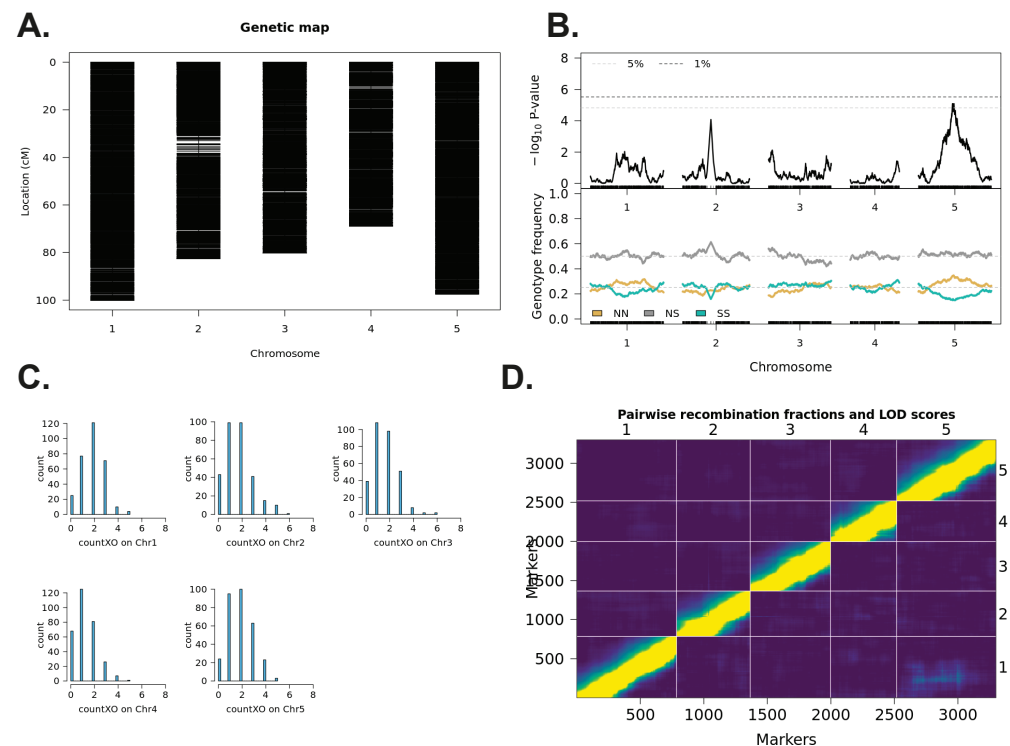
	Gains (%)		Losses (%)	
Line	NN×SS	SS×NN	NN×SS	SS×NN
N	0.03	0.04	0.30	0.11
S	0.06	0.03	0.27	0.23

S6 Table. Epimutation rates using data from S14 Fig.

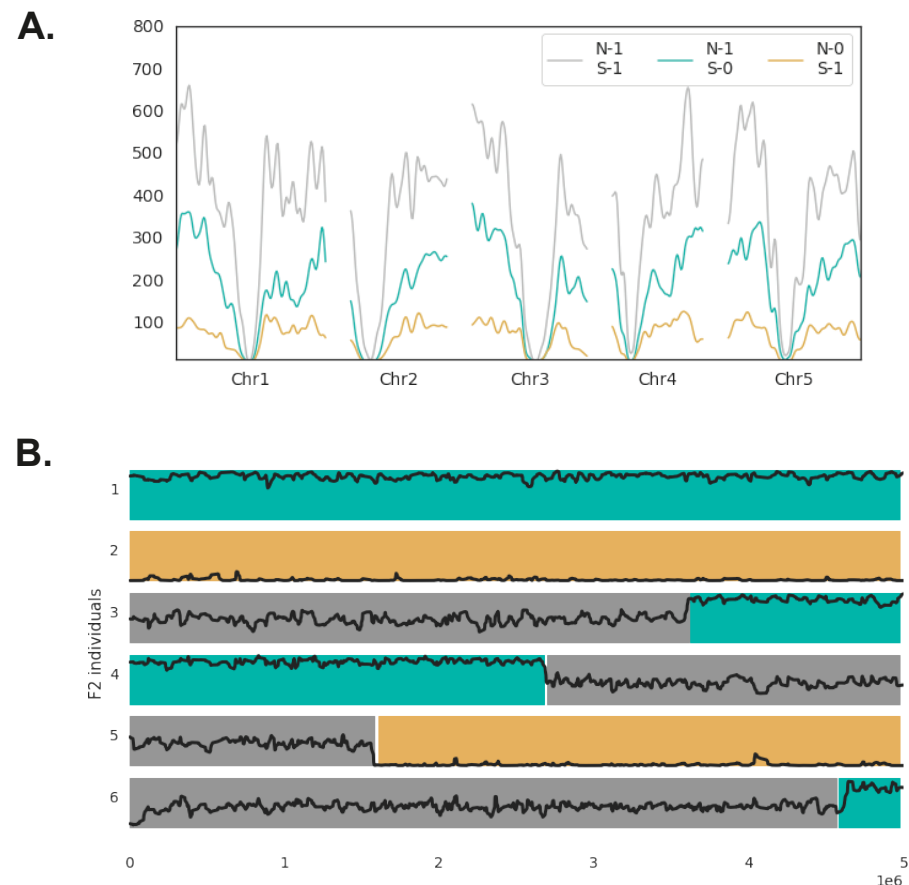
	Gains (%)		Losses (%)	
Line	NN×SS	SS×NN	NN×SS	SS×NN
N	0.55	0.38	0.71	0.20
S	0.40	0.27	0.68	0.61



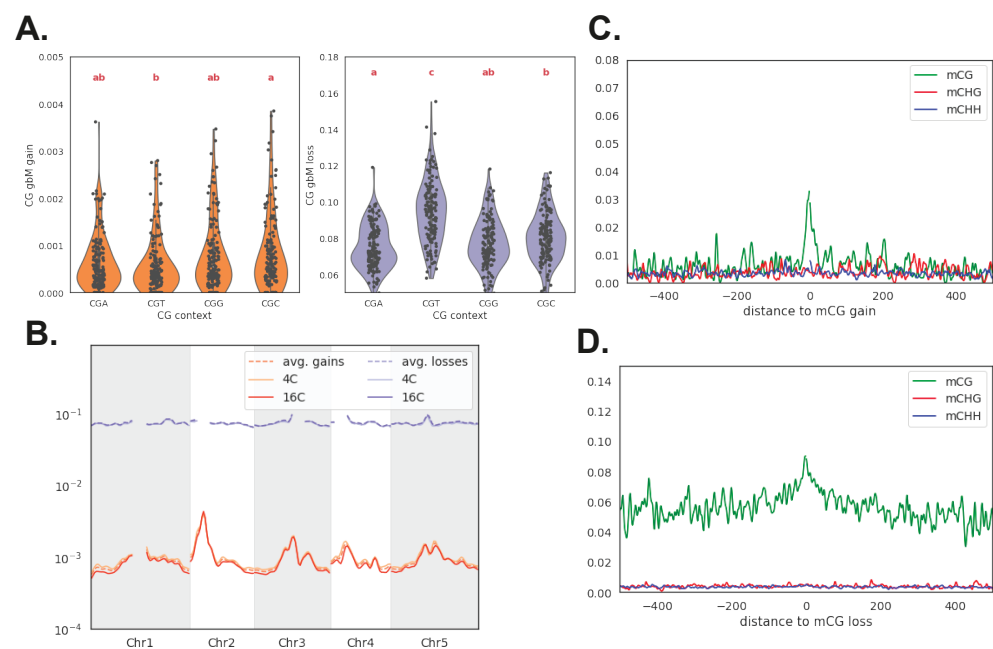
S1 Fig. Experimental design and residual heterozygosity. (A) Reciprocal F2 design. **(B)** The left panel shows evidence for residual heterozygosity in the parental lines in the 1001 Genomes data. The right panel shows region where different SNPs are segregating in the reciprocal F2 populations.



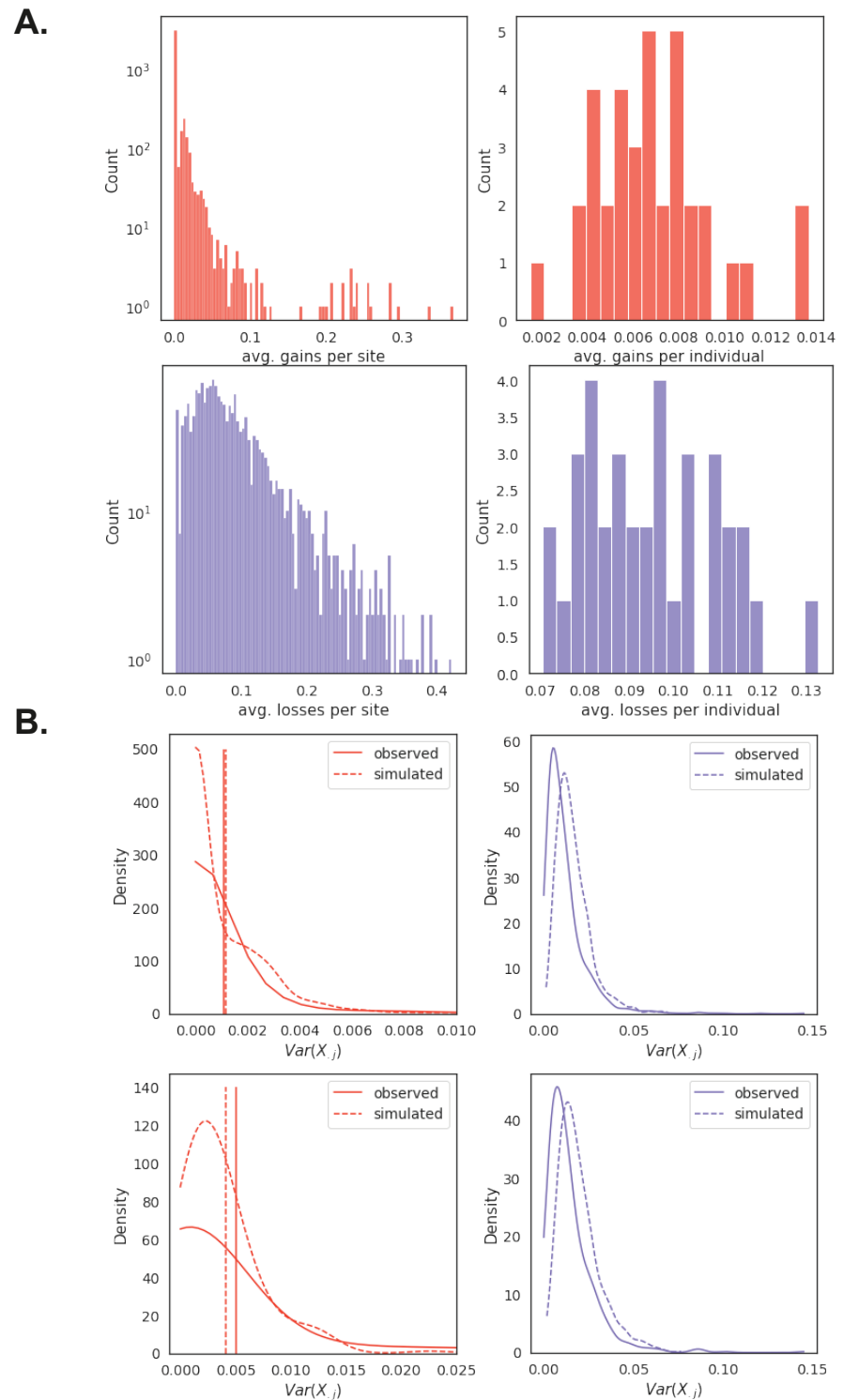
S2 Fig. Genetic map construction for the NN \times SS cross. (A) Genetic map and markers. **(B)** Segregation distortion in the cross and genotype frequencies across the chromosome. **(C)** Number of crossovers per chromosome in the genetic map. **(D)** Pairwise recombination fraction (upper left triangle) and LOD scores for the markers.



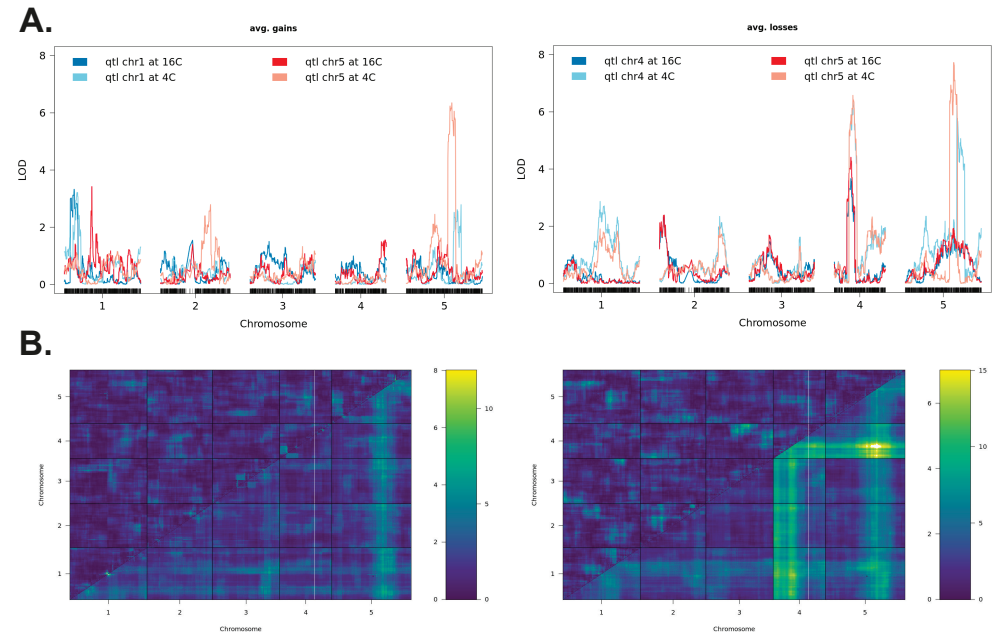
S3 Fig. Mendelian segregation for gene-body methylation. (A) Distribution of methylated CG sites the genome in 200 kb windows, separately for 213178 sites methylated in both parents (N-1 S-1), 109868 sites methylated only in the northern parent (N-1 S-0), and 39682 sites methylated only in the southern parent (N-0 S-1). (B) Genotype and relative methylation levels for 6 F2 individuals along chromosome 1. Genotypes are given by colors (NN is turquoise; SS is yellow; NS is grey), relative methylation levels by black curve.



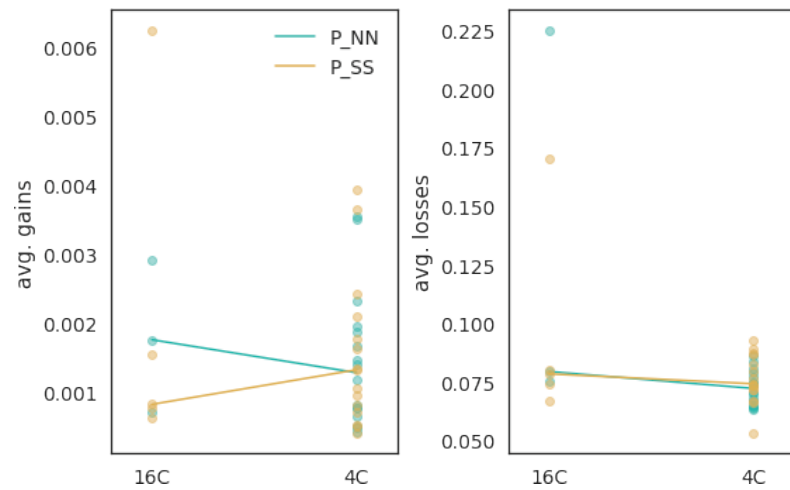
S4 Fig. Patterns of somatic deviations. (A) Gains and losses separated by four contexts (CGA, CGT, CGG and CGC). (B) Average gains and losses across the genome at different temperatures. (C) Methylation levels at (gain at previously unmethylated) CG, CHG and CHH sites near a mCG gain site. (D) Methylation levels at CG (loss at previously methylated), CHG and CHH sites near a CG loss site.



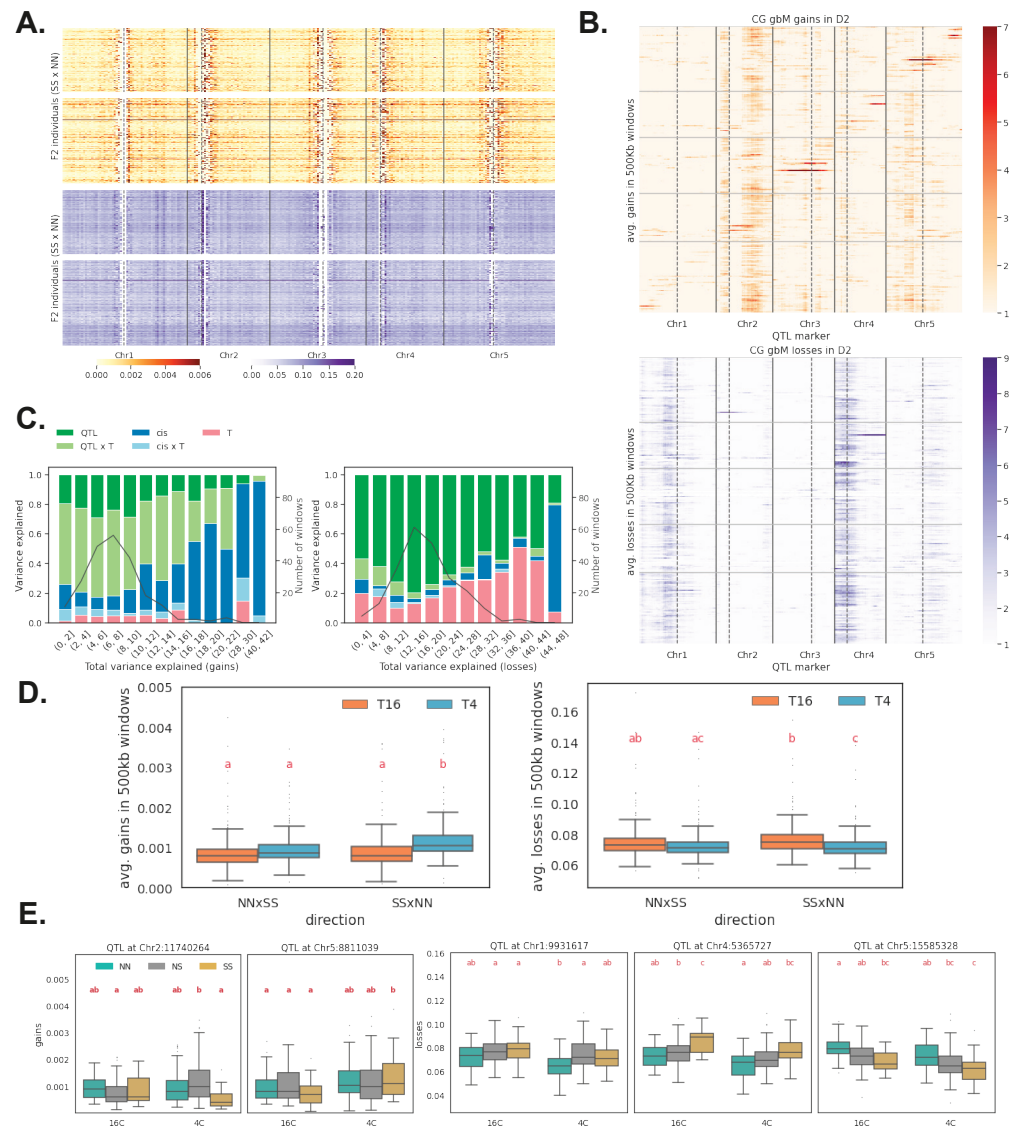
S5 Fig. Modeling somatic deviations. (A) Distribution of average deviations per site and per individual using data from chromosome 1 NN genotypes as an example. (B) Distribution of the variance between individuals across sites, $\text{Var}(X_{.j})$, in data and in simulations. Top row shows the distribution for all sites, bottom row only for sites that are differentially methylated between N and S.



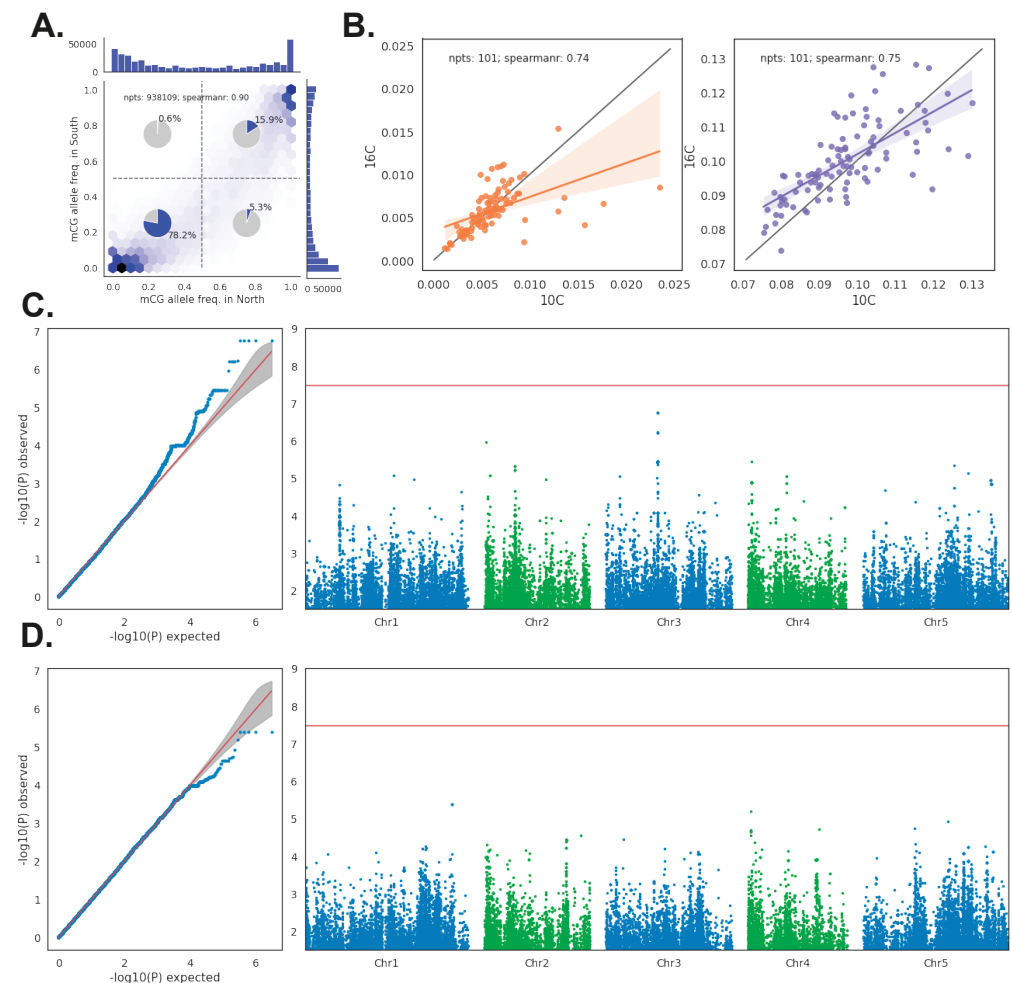
S6 Fig. Composite Interval Mapping for average deviations. (A) Composite Interval Mapping was applied to four different gain phenotypes and four different loss phenotypes in order to refine peaks. For each of the four major QTL identified by combining results across 500 kb windows (two for gains and two for losses, see Fig 4) deviations were averaged over regions showing QTL effect at two temperatures. (B) Testing for epistasis on the QTLs for somatic deviations (using the “scantwo” function in R/qtl). Two QTLs on Chr1 and Chr5 for gains and three QTLs on Chr1, Chr4 and Chr5 for losses. The bottom triangle is the LOD scores for the full model including the interaction effect, upper triangle is LOD scores for only the interaction.



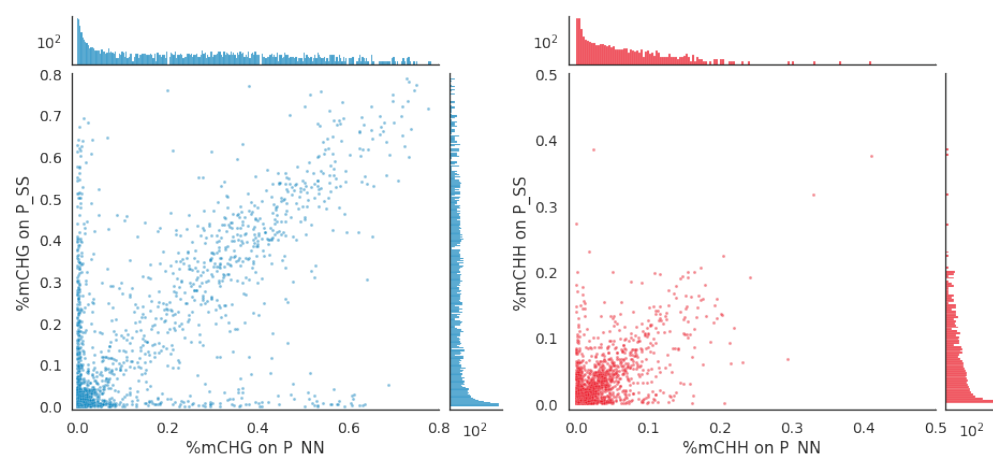
S7 Fig. Deviations in parental strains. Reaction norms for average gains and losses for parental strains.



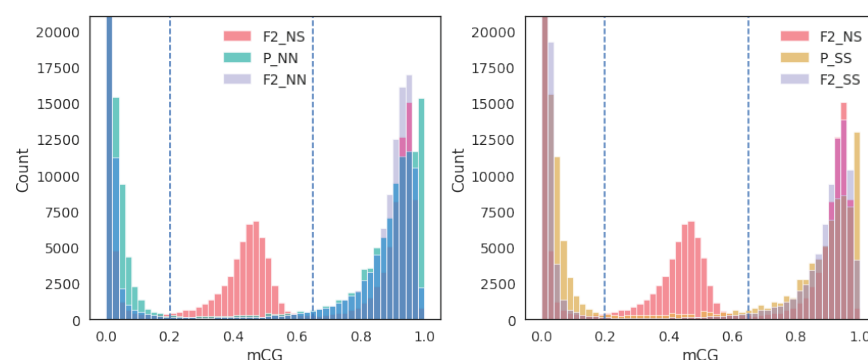
S8 Fig. Genetic architecture of deviations in reciprocal cross (SSxNN). (A) Average deviations in 500 kb windows across genome (cf. Fig 3). (B) QTL mapping for gains and losses (cf. Fig 4). (C) Variance-partitioning results (cf. Fig 4). (D) Temperature effects on average gains and losses (in NN background) for both directions. (E) Genotypic effects for two gain QTL and three loss QTLs (cf. Fig 5).



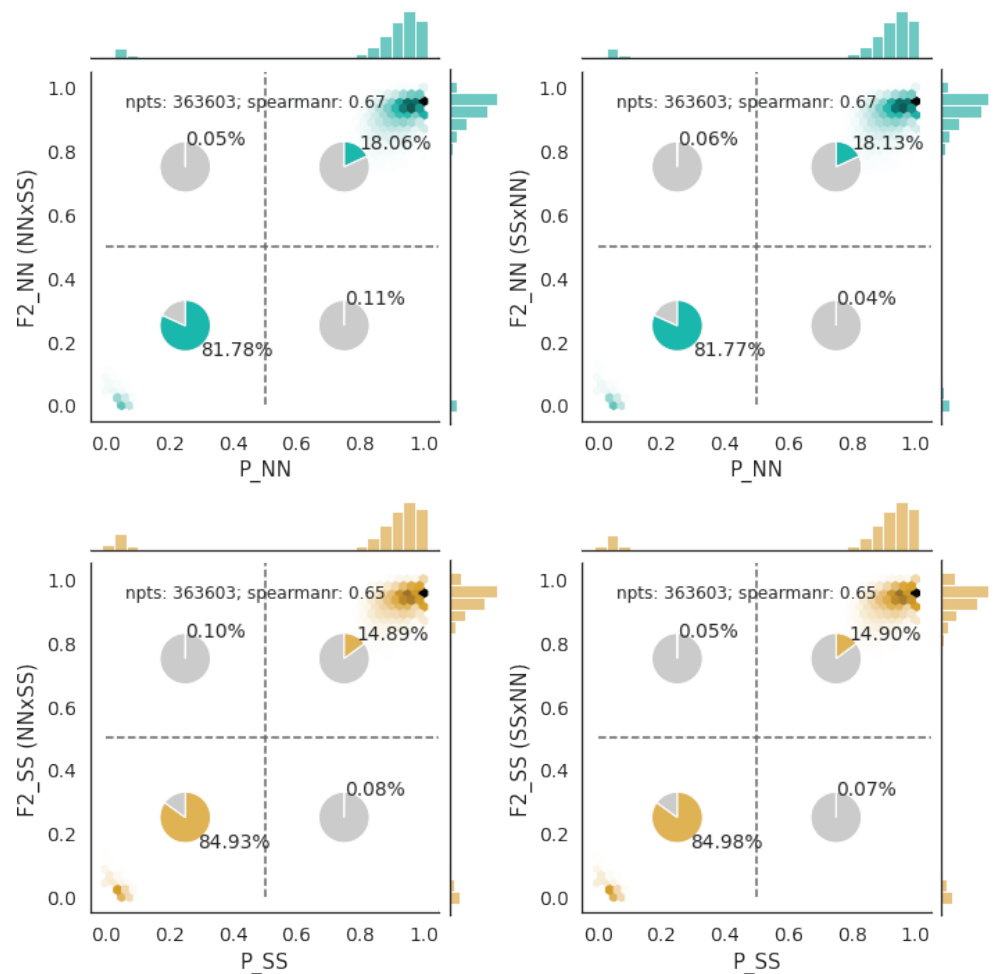
S9 Fig. GWAS of deviations in Swedish *A. thaliana*. (A) mCG allele frequencies in populations from northern and southern Sweden [3]. (B) Correlation between genome-wide deviations between 10°C and 16°C. (C) GWAS for genome-wide gains. (D) GWAS for genome-wide losses.



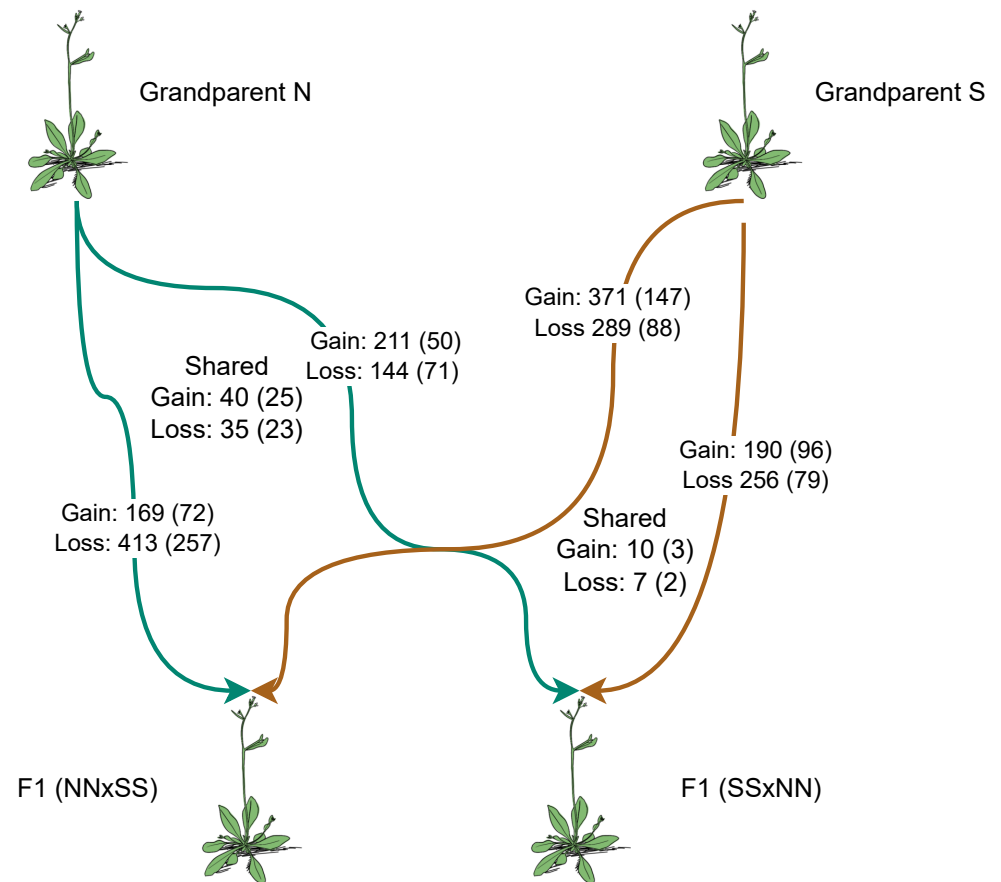
S10 Fig. Gene-body methylation %mCHG and %mCHH on annotated protein coding genes (Araport 11) in parental lines N and S. We filtered out genes having any non-CG methylation on the gene-bodies to determine the gbM genes.



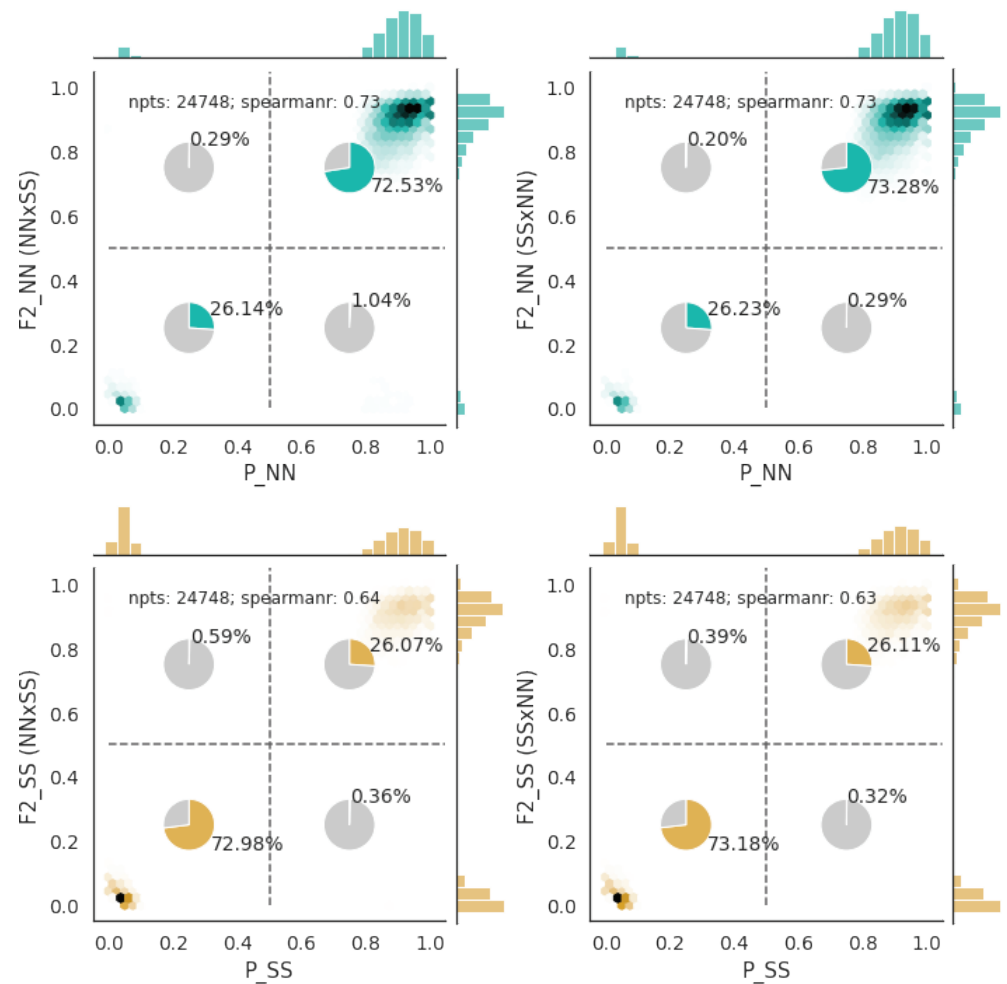
S11 Fig. Possibly heterozygous mCG sites in grandparents Histograms for methylation levels on gbM sites averaged across parental values (N in the left panel and S on the right panel), F2 individuals homozygous for the same ancestry, and F2 individuals heterozygous for ancestry. There are far more sites with intermediate values in the parental than in the homozygous F2 data, although the former is also supposed to be homozygous.



S12 Fig. Trans-generational epimutations. The plots compare average gbM for parents with average gbM for F2 individuals homozygous for the parental ancestry across sites, separately for the two cross-directions. Only data from chromosome 5 was used as all other chromosomes showed evidence of residual heterozygosity in the southern parental line (S1 Fig).



S13 Fig. Trans-generational epimutations along lines of descent. The trans-generational epimutation from S12 Fig are shown for each line-of-descent in the cross. "Shared" refers to the number of changed sites that are shared between the directions of the cross, separately for the northern and southern ancestry. The numbers in parentheses are for the sites that are differentially methylated sites between N and S (S14 Fig).



S14 Fig. Transgenerational epimutations. Same plots as S12 Fig but only sites that differ in methylation between the parental lines were used.