

1 context, and thus the 'epigenetic potential' of a population evolves at the sequence level (Kilvitis
2 et al. 2017, see also Adrian-Kalchhauser et al. 2020).

3 It is also established that the epigenetic conformation of the genome affects the propensity for
4 sequence change. Notably, DNAm can influence mutation rates due to the higher susceptibility
5 of methylated Cs to spontaneous deamination to form thymine (Xia et al. 2012; Poulos et al.
6 2017; Zhou et al. 2020). In mammals CpG mutation rates have been estimated at 10-50x higher
7 than in other sequence contexts (Walser and Furano 2010). CpG mutation rate also has a
8 nuanced relationship with methylation levels, as CpG sites with the highest mutation rates in
9 human populations were observed to have low-to-intermediate methylation levels in cultured
10 cells (Xia et al 2012). Therefore, by influencing sequence evolution, epigenetic variation may
11 have unappreciated roles in the emergence of genomic novelties and adaptations (Storz et al.
12 2019; Guerrero-Bosagna 2020), as well as mediating environmental influences on sequence
13 evolution (Guerrero-Bosagna 2012; Lu et al. 2021).

14 Despite the interdependence of DNAm and sequence variation, the potential importance of this
15 link in local adaptation has been largely overlooked. For example, typical workflows for
16 detection of differential methylation tend to exclude CpG sites that are not detected at a certain
17 coverage in a certain proportion of individuals (e.g. Akalin et al. 2012), and therefore it could be
18 assumed that genetic diversity of those sites is irrelevant. However, these sites may
19 nevertheless harbour genetic variants in the population, the relative frequencies of which may
20 be informative about evolutionary forces acting on methylation state, potentially allowing further
21 dissection of the manner in which DM evolves in the context of local adaptation. Indeed,
22 methylation sites within certain promoters have already been shown to exhibit selective sweep
23 signatures in *Arabidopsis* (Shirai et al. 2021). Epigenetic diversification is one of many possible
24 routes to local adaptation (e.g. Smith et al. 2016) but may occur in conjunction with others, such
25 as selection on discrete new mutations (hard sweeps) or on standing genetic variation (soft
26 sweeps) (Bernatchez 2016; Hermisson and Pennings 2017). For example, if epigenetic
27 modifications at multiple loci could confer similar adaptive benefit, epigenetic diversification
28 could occur in conjunction with a soft sweep. Furthermore, given the heightened mutation rate
29 of methylated Cs and its complex relationship with methylation levels (Xia et al, 2012),
30 acquisition of methylation in a divergent population or a change in methylation level may
31 influence mutation rates at affected sites. There is therefore a need to examine the relationships
32 between differential methylation and nucleotide diversity in divergent populations.

1 typically lower in freshwater populations, likely due to a past bottlenecks (Terekhanova et al.
2 2014), patterns of nucleotide diversity at methylation sites, and DMCs specifically, have not
3 been addressed. Shifts in nucleotide diversity at DMCs may be informative about the
4 evolutionary forces acting on DNAm during local adaptation, namely the tightening or
5 loosening of selective constraint on methylation state in a new environment.

6 Here, we combine epigenetic and genetic data from the White Sea stickleback population
7 complex to study the interactions between methylation differences and nucleotide diversity
8 during freshwater colonisation. We examined nucleotide diversity in relation to methylation
9 divergence, variance, and the environmental inducibility of methylation state, considering both
10 variance and inducibility as indicators of the relative stringency DNAm regulation.

11 **Results**

12 *Elevated nucleotide diversity accompanies differential methylation but level depends on the*
13 *direction of methylation change*

14 For generation of DNAm and genome sequence data, respectively, both Artemov et al. and
15 Terekhanova et al. sampled freshwater fish from the same lake (Lake Mashinoye) and marine
16 fish from nearby coastal locations in the Kandalaksha gulf. A combined analysis of samples
17 from these two datasets therefore allowed us to identify differentially methylated cytosines in
18 CpG context (DMCs) between marine and freshwater populations and examine the nucleotide
19 diversity of these sites in separate samples of those populations (**Fig. 1**). The RRBS data
20 (Artemov et al. 2017) derived from gill tissue from a total of 11 individuals. These included six
21 individuals as part of the main population comparison (N = 3 per population) and a further five
22 experimental treatment animals that were used for a subsequent analysis of site inducibility.
23 After filtering to remove sites with C-T/G-A SNPs detected in RRBS individuals, which could
24 otherwise lead to spurious counts of unmethylated Cs, the analysis included just over 1 million
25 CpG sites with at least 5x alignment coverage in all individuals, comprising approx. 6.9% of
26 CpGs in the genome. The pool-seq data (Terekhanova et al. 2014, 2019) comprised sequenced
27 material of two pools containing 12 marine and 10 freshwater individuals, and with genome
28 coverage (high quality alignments) of 10.4x and 8x, respectively.

29 We first derived measures of nucleotide diversity of differentially methylated (DMCs) and non-
30 differentially methylated DNAm sites (Non-DMCs) in the form of π (average number of pairwise
31 differences within population) (Nei and Li 1979), Watterson's θ (population-scaled mutation rate)
32 (Watterson 1975), and Tajima's D (Tajima 1989). π and θ are complementary measures of

1 These patterns of elevated nucleotide diversity were not driven by enrichment for DMCs in
2 regions of high diversity. Rather, elevated π of DMCs was found to be strongly localised around
3 individual DMCs (**Fig. S3**). The pattern was largely consistent across different genomic features
4 including CpG islands, gene bodies, promoters, and intergenic regions (**Fig. S3**). No clear
5 pattern of localised elevated diversity was observed for DMCs which fell within differentially
6 methylated regions (DMRs) (**Fig. S3**), however only a fraction of FW-hypo (approx. 17%) and
7 FW-hyper DMCs (approx. 12%) fell within DMRs.

8 Next, we determined which mutation type(s) were most likely to be driving the elevated π of
9 DMCs, and specifically whether this was driven by an over-abundance of C-T transitions. To this
10 end, the percentages of sites in each category harbouring biallelic SNPs of different types (C-
11 T/G-A, C-A/G-T, or C-G/G-C) were calculated from the pool-seq data. The majority of SNPs
12 were C-T/G-A, comprising 90% of SNPs across all categories in marine and 94% in freshwater.
13 The proportion of sites harbouring biallelic C-T/G-A SNPs across the three categories of
14 methylation site and two populations showed a similar pattern to that of π , with FW-hyper sites
15 harbouring the highest proportion of C-T/G-A SNPs in both populations (**Fig. 2D**). Marine had
16 more C-T/G-A SNPs than freshwater in the Non-DMC (paired Wilcoxon test, $p < 0.001$) and FW-
17 hyper categories ($p < 0.001$) but not the FW-hypo category, in which FW and marine had similar
18 proportions of C-T/G-A SNPs. Meanwhile, the percentage of other SNP types showed no clear
19 differences between the site categories. Therefore, the generally increased nucleotide diversity
20 amongst DMCs relative to Non-DMCs seemed to be driven by a greater occurrence of C-T
21 mutations.

22 *Higher nucleotide diversity of infrequently-methylated DMCs*

23 The finding that sites which gained methylation in freshwater (FW-hyper) had the highest π and
24 highest proportion of C-T/G-A mutations in marine (**Fig. 2**) was contrary to expectations, as
25 these sites would be expected to be infrequently methylated in marine and therefore not at high
26 risk of mutation via deamination. We therefore tested the relationship between π and the
27 distributions of mean percentage of methylation (hereafter Mean{PM}) across the three site
28 categories. Here, Mean{PM} refers to the average percentage of copies on which the C is
29 methylated, or in other words the average frequency of the methylation mark among copies of
30 given CpG site. We found that Non-DMCs displayed a bimodal density distribution, with most
31 sites either very frequently (>75%) or very infrequently methylated (**Fig. 3A**, left). Meanwhile,
32 the distributions of Mean{PM} of DMCs were markedly different to those of Non-DMCs. FW-
33 hypo sites were characterised by a shift from mostly high Mean{PM} in marine to mostly

1 intermediate Mean{PM} in FW (**Fig. 3A**, middle). Mirroring this pattern, FW-hyper sites were
2 characterised by a shift from low-intermediate Mean{PM} in marine to high Mean{PM} in FW
3 (**Fig. 3A**, right).

4 Because pool-seq data are not appropriate for estimating π at the level of a single site, we used
5 a ranking procedure to examine the relationship between Mean{PM} and π . Sites were divided
6 into ranks according to their Mean{PM}, with higher ranks containing sites with higher
7 Mean{PM}. This ranking was performed separately for each population and each site category,
8 and a measure of π obtained for each rank. The relationship between the rank-level Mean{PM}
9 and π was clearly non-monotonic for Non-DMCs, with the highest values appearing at low to
10 intermediate Mean{PM} of around 25% (**Fig. 3B**, left). The higher π of FW-hypo sites in
11 Freshwater appeared to be driven largely by sites in the low-intermediate range (**Fig. 3B**,
12 middle). Amongst FW-hyper sites, those with the highest Mean{PM} clearly contributed to the
13 lower π of these sites in FW (**Fig. 3B**, right).

14 We also examined the relationship of π with the population difference in Mean{PM} (i.e. the
15 extent of hypo- or hypermethylation). We observed that among FW-hypo sites, the freshwater
16 population had the largest increases in the π where the hypomethylation was strongest (**Fig. 3C**
17 and **Fig. 3D**). FW-hyper sites also increased in π with the extent of hypermethylation (**Fig. 3C**),
18 but this also corresponded with greater loss of π in FW (**Fig. 3D**). Meanwhile, sites with larger
19 difference in Mean{PM} in either direction (methylation loss or gain) had higher F_{st} (**Fig. 3E**).

20 *High nucleotide diversity accompanies high variability in ancestral methylation*

21 Considering that sites with intermediate Mean{PM} are liable to have more variable methylation
22 frequency than those with very low or very high Mean{PM}, we also considered the relationship
23 between π and the standard deviation of percentage methylation (hereafter SD_{meth}). We
24 predicted that sites with more variable methylation would have higher nucleotide diversity,
25 reasoning that the methylation state of these sites is not stringently controlled and therefore
26 mutations at these sites may have little impact on function. We first examined the distributions of
27 SD_{meth} values of sites in the Non-DMC, FW-hypo and FW-hyper categories. Non-DMCs were
28 largely invariable, with slightly more variable methylation in freshwater compared to marine (**Fig.**
29 **4A**, left), consistent with the observations of Artemov et al. (2017). FW-hypo sites were
30 characterised by a pronounced increase in SD_{meth} from ancestral to derived population, shifting
31 from a left-skewed distribution in marine to a Gaussian-like distribution in freshwater (**Fig. 4A**,

1 Hu and Barrett 2022). As such, we considered environmental inducibility in a different context, in
2 that the degree of environmental inducibility of methylation state is (inversely) indicative of the
3 degree of intrinsic regulation. We therefore use the term ‘inducibility’ loosely to refer to the
4 sensitivity of a site to methylation change in response to the environment, regardless of its
5 potential adaptive importance. We found that elevated π of and Fst of DMCs was driven by sites
6 that were environmentally inducible (**Fig. 5A, B, E**), further supporting a hypothesis of relaxed
7 regulation and relaxed selective constraint at sites that are responsive to environmental
8 conditions. Furthermore, the increased π among FW-hypo sites in FW relative to marine was
9 driven by sites that were induced only in FW, i.e. those not induced in the ancestral population,
10 suggesting that nucleotide diversity is more likely to accumulate at sites where intrinsic control
11 of methylation is relaxed (and therefore more sensitive to the environment). Indeed, the positive
12 correlation between π and the degree of inducibility (**Fig. 5C**) suggests that the more sensitive
13 the methylation state is to the environment, the more likely mutations are to be selectively
14 neutral. Therefore, shifts in inducibility (in addition to shifts in methylation variance, as discussed
15 above) may precede shifts in nucleotide diversity. Our results suggest that the majority of
16 environmentally inducible sites are simply ‘blowing in the wind’ and do not have important
17 functions for plasticity which would constrain nucleotide diversity. Nevertheless, in their analysis
18 of Baltic Sea sticklebacks, Heckwolf et al. (2020) observed that the Fst of induced sites (marine
19 fish responsive to lower salinity) depended on the direction of the induced change. Sites that
20 were induced to the ‘evolved’ methylation state observed in the derived freshwater population
21 had lower Fst than those that were induced in the opposite direction. This suggests that some
22 environmentally inducible sites are indeed constrained by selection due to the importance of site
23 plasticity. Here, we did not consider the direction of inducible change, merely considering
24 inducibility as a proxy for the relative weakness of intrinsic regulation.

25 Again, these observations could also be reconciled with the scenario of a soft sweep, as it is
26 also possible that plasticity of only some methylation sites is necessary to confer adaptation. In
27 other words, plasticity of multiple sites provides multiple alternate routes to adaptation. As many
28 methylation sites would therefore be redundant, they could be lost to mutation without
29 detrimentally affecting the organism’s capacity for adaptive plasticity.

30 *Limitations and future directions*

31 Our analyses have revealed striking associations between genetic and epigenetic variation in
32 divergent stickleback populations. However, we must acknowledge limitations including the

1 unmethylated Cs that were bisulfite-converted to Ts. Although we used a combination of three
2 SNP-callers designed for BS-seq data (see methods), we cannot be certain that some
3 differential methylation was not the result of SNPs that these algorithms failed to detect (see
4 Lindner et al. 2022).

5 In a broader context, our study is limited in that we only examined one population pair. It is
6 therefore currently not known whether the patterns we observed occur more broadly across
7 different local adaptations (in stickleback and other species) or whether they are idiosyncratic to
8 the relatively recent colonisation event considered in this study (~700 years). The existence of
9 far older populations, such as those in the Japanese archipelago which are estimated to have
10 colonised freshwater ~170,000 years ago (Kakioka et al. 2020), raises the question as to the
11 fate of differential methylation over longer periods. Over time, for example, the initially
12 heightened methylation variance may return to a less variable state due to refinement of
13 methylation states via selection or the removal of the CpG sites via accumulation of C-T
14 transitions. Alternatively, no substantial accumulation of mutations over time would suggest that
15 the heightened diversity of FW-hypo sites reflects standing genetic variation.

16 If, indeed, heightened methylation variance arises due to relaxed control of methylation state,
17 the mechanisms by which this could occur are not known. Artemov et al. (2017) suggested that
18 mutations in genes encoding epigenetic regulators may underlie increased methylation
19 variance, but did not identify any known epigenetic regulators in the vicinity of genomic regions
20 differentiating marine and freshwater populations in the White Sea region. *Trans*- and *cis*-
21 meQTL have however been identified in stickleback (Hu et al. 2021), some of which are indeed
22 in the vicinity of genomic regions of high F_{st} between marine and freshwater populations.
23 Differential selection on *trans*-meQTL in particular could have knock on effects on methylation
24 sites across the genome.

25 While our study considered only genetic variation in the form of SNPs at CpG sites themselves,
26 DNAm is associated with other types of genome sequence alterations. These include
27 mutations in non-CpG context (Walser and Furano 2010), recombination rate variation (Mirouze
28 et al. 2012), and structural variation including copy number and transposable element variation
29 (Guerrero-Bosagna 2020). Indeed, the role of structural variation in local adaptation is
30 increasingly appreciated and major inversions, transposable elements and copy number
31 variants are all proposed to have played a role in stickleback freshwater adaptation (Reid et al.
32 2021). The interplay between epigenetic variation and other forms of genetic variation therefore

1 run accession: SRR869609), while marine fish were collected from the Kandalaksha gulf as part
2 of the 2014 study and a subsequent 2019 study (Terekhanova et al. 2019). We selected the
3 'White Sea, WSBS' sample from Terekhanova et al. (2019) (SRR7470095) as the marine
4 sample for our comparison, given that it has a similar pool size to the Mashinnoye sample (12
5 vs 10) and a similar number of 100bp paired reads (64,176,648 vs 62,016,859 after quality
6 trimming). Sequence files were obtained in FASTQ format from the Sequence Read Archive
7 (SRA) and European Nucleotide Archive (ENA).

8 *Data processing: RRBS*

9 Raw RRBS reads were trimmed using TrimGalore v0.6.6 using default settings. Alignment to
10 the Three-spined stickleback v.5 assembly (Nath et al. 2021) and subsequent methylation
11 calling were carried out using Bismark v0.22.3 (Krueger and Andrews 2011) with Bowtie2
12 v2.3.4.1 as the aligner (Langmead and Salzberg 2012). Methylation calls were not strand-
13 specific. To remove sites harbouring C-T/G-A SNPs which otherwise contribute erroneous
14 counts of non-methylated Cs, we ran three SNP-callers on each sample: BS-SNPper v1.1 (Gao
15 et al. 2015), Biscuit v0.3.14 (<https://github.com/zhou-lab/biscuit>), and CGmap-tools v0.1.2 (Guo
16 et al. 2018). We then compiled the coordinates of all sites harbouring C-T/G-A SNPs detected in
17 any of the individuals by any of the SNP-callers (either homo- or heterozygous), and removed
18 these sites from the Bismark coverage files containing the methylation counts (counts of Cs and
19 Ts at each position). This approach detected 75% of C-T/G-A that were detected at high
20 frequency in the freshwater pool-seq sample (Fig. S5). Further details of SNP calling from
21 RRBS are provided in the supplementary methods.

22 *Data processing: Pool-seq*

23 Raw reads were trimmed with Trimmomatic v0.36 (Bolger et al. 2014) with the option
24 SLIDINGWINDOW:4:20 and otherwise default parameters. Only reads which remained paired
25 after trimming were kept. Reads were mapped to the Three-spined stickleback v.5 assembly
26 with Bowtie2 v2.3.4.1 with default parameters (Langmead and Salzberg 2012). Sambamba
27 v0.7.1 (Tarasov et al. 2015) was used to filter the alignments to retain those with MAPQ \geq 20
28 and to remove PCR duplicates. This resulted in 46,588,899 and 35,691,831 high quality
29 alignments from marine and FW samples, equating to average genome coverage of 10.4x and
30 8x, respectively. Samtools v0.1.18 (Danecek et al. 2021) was used to generate a pileup file from
31 each BAM file, as required for the Popoolation and Popoolation2 toolkits.

32 *Identification of differentially methylated CpG sites (DMCs) and subsampling of Non-DMCs*

1 per-chromosome basis was because Popoolation's estimates of π are accurate over large
2 numbers of sites, but not at the single site level (Kofler, Orozco-terWengel, et al. 2011). Each
3 site was labelled with its chromosome and its category within the analysis (e.g. chr1 FW-hypo)
4 and the labelled category was entered as the 'gene ID' in a GTF file, such that variance-at-
5 position.pl, which was developed to calculate diversity statistics per-gene, was instructed to
6 calculate π for each combination of chromosome and site category. A similar procedure was
7 used to obtain π for sites ranked according to (difference in) mean percentage methylation
8 (Mean{PM}), (difference in) SD of percentage methylation (SD_{meth}), and absolute inducibility,
9 whereby ranks were assigned using the bin() function from the OneR package, specifying 50
10 ranks each time. Sites were then labelled in the GTF according to their rank (regardless of
11 chromosome), such that a single value of π was obtained for each rank. Variance-at-position.pl
12 from Popoolation was run with the parameters --min-qual 20 --min-coverage 3 --min-count 2.
13 The majority of sites met the requisite 3x coverage for inclusion in nucleotide diversity estimates
14 of marine (99%) and FW (93%). For the analysis of nucleotide diversity as a function of absolute
15 inducibility, one rank was excluded from the FW population due to insufficient coverage (<60%
16 of sites with 3x coverage).

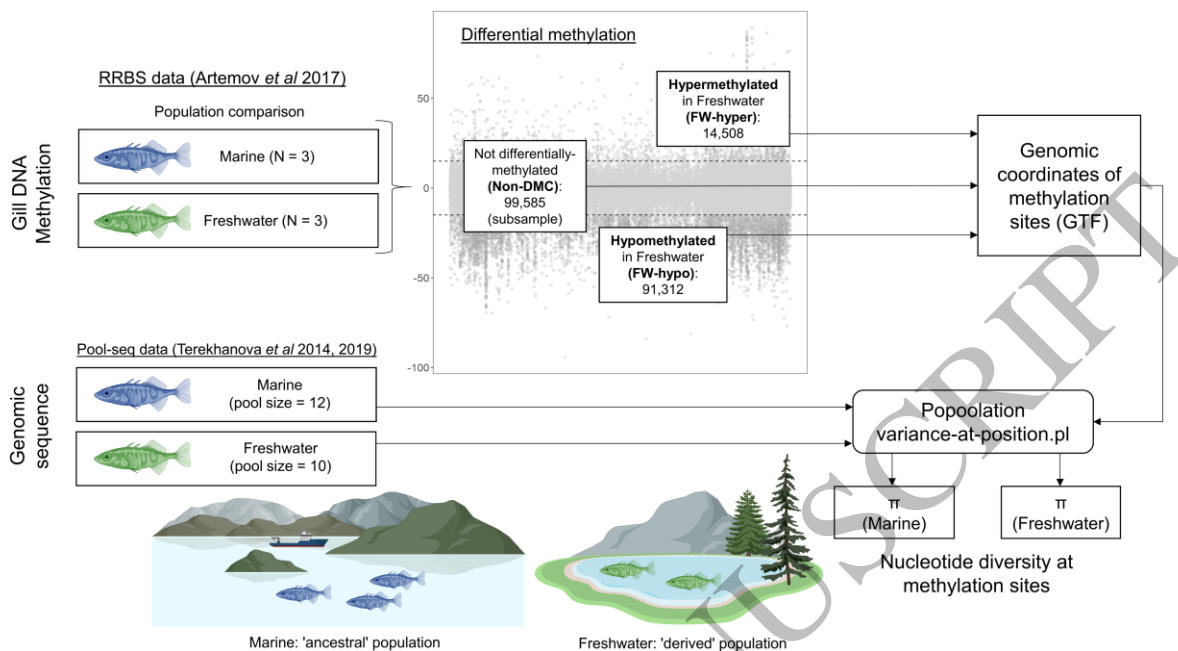
17 Fst for different categories of methylation sites (including ranked sites) were obtained using the
18 Popoolation2 toolkit (Kofler, Pandey, et al. 2011). 'Gene-wise' .sync files were obtained from the
19 pileup files using coordinates in the abovementioned gtf files and were used as input for the 'fst-
20 sliding.pl' script which was run with parameters --min-count 2 --min-coverage 3 --pool-size 22 --
21 min-covered-fraction 0 --max-coverage 1000 --window-size 1000000 --step-size 1000000.

22 *Percentage of sites with SNPs*

23 To obtain the % of sites within each result category (Non-DMC, FW-hypo, and FW-hyper)
24 harbouring SNPs of different types (C-T/G-A and non-C-T/G-A) in the pool-seq samples, we
25 filtered the BAM files of each population to retain alignments corresponding with the positions of
26 interest. We then ran GATK HaplotypeCaller (McKenna et al. 2010) with the --sample-ploidy set
27 to the pool size x 2 (24 for marine and 20 for freshwater), and otherwise default settings. The
28 subsequent VCF file was then filtered using bcftools v1.10 to retain only biallelic SNPs at the
29 sites of interest. We subsequently extracted from the VCF a list of reference and alternate
30 alleles at sites of interest harbouring biallelic SNPs. We were therefore able to assign SNPs as
31 either 'C-T/G-A' or 'other', and calculate the % of sites in each result category harbouring
32 biallelic SNPs of one of those two classes.

1 **Fig. 5. Nucleotide diversity of differentially methylated sites in relation to their capacity**
2 **for induced methylation change in response to environmental salinity. (A)** Additional
3 RRBS data deriving from experimental salinity treatments performed by Artemov et al (2017)
4 (marine fish placed in freshwater and freshwater fish placed in saltwater) were used to identify
5 sites that were inducible in response to salinity change in either of the two populations. **(B)** % of
6 sites in the Non-DMC, FW-hypo, and FW-hyper categories that were induced in response to
7 salinity change in the marine (blue) and freshwater (green) populations. *P*-values derived from
8 paired Wilcoxon tests. **(C)** Per-chromosome estimates of π for FW-hypo and FW-hyper sites
9 divided according their capacity for induced gill methylation change in response to a change in
10 environmental salinity, considering sites that were induced in neither of the populations, either of
11 the two populations, or only in one of the two populations (marine or freshwater). π of Non-
12 DMCs is shown in separate panel for comparison. *P*-values derived from paired *t*-tests. **(D)** π of
13 inducible DMCs that were ranked according to their mean absolute induced change in
14 percentage of methylation (i.e. regardless of the direction). Within each population, only sites
15 that were significantly differentially methylated in response to salinity (mean difference in
16 percentage of methylation ≥ 15 or ≤ -15 , $p < 0.05$) were considered. Separate ranks were
17 obtained for marine and FW and a single π estimate obtained for each rank. **(E)** Per-
18 chromosome estimates of pairwise *F*_{st} (freshwater vs. marine) of FW-hypo and FW-hyper sites
19 divided according to their capacity for induced methylation change. *F*_{st} of Non-DMCs is shown
20 in separate panel for comparison. *P*-values derived from paired *t*-tests. **(F)** Pairwise *F*_{st} of
21 inducible DMCs that were ranked according to their mean absolute induced change in
22 percentage of methylation (i.e. regardless of the direction of the change). For **(D)** & **(F)**, 50 ranks
23 were used for marine and 49 for freshwater. Each rank contains an average of 227 sites for
24 Marine and 228 sites for Freshwater. Trend lines derived from linear models and ribbons show
25 SEM.

26
27



1
2
3

Figure 1
159x91 mm (.97 x DPI)

ACCEPTED MANUSCRIPT

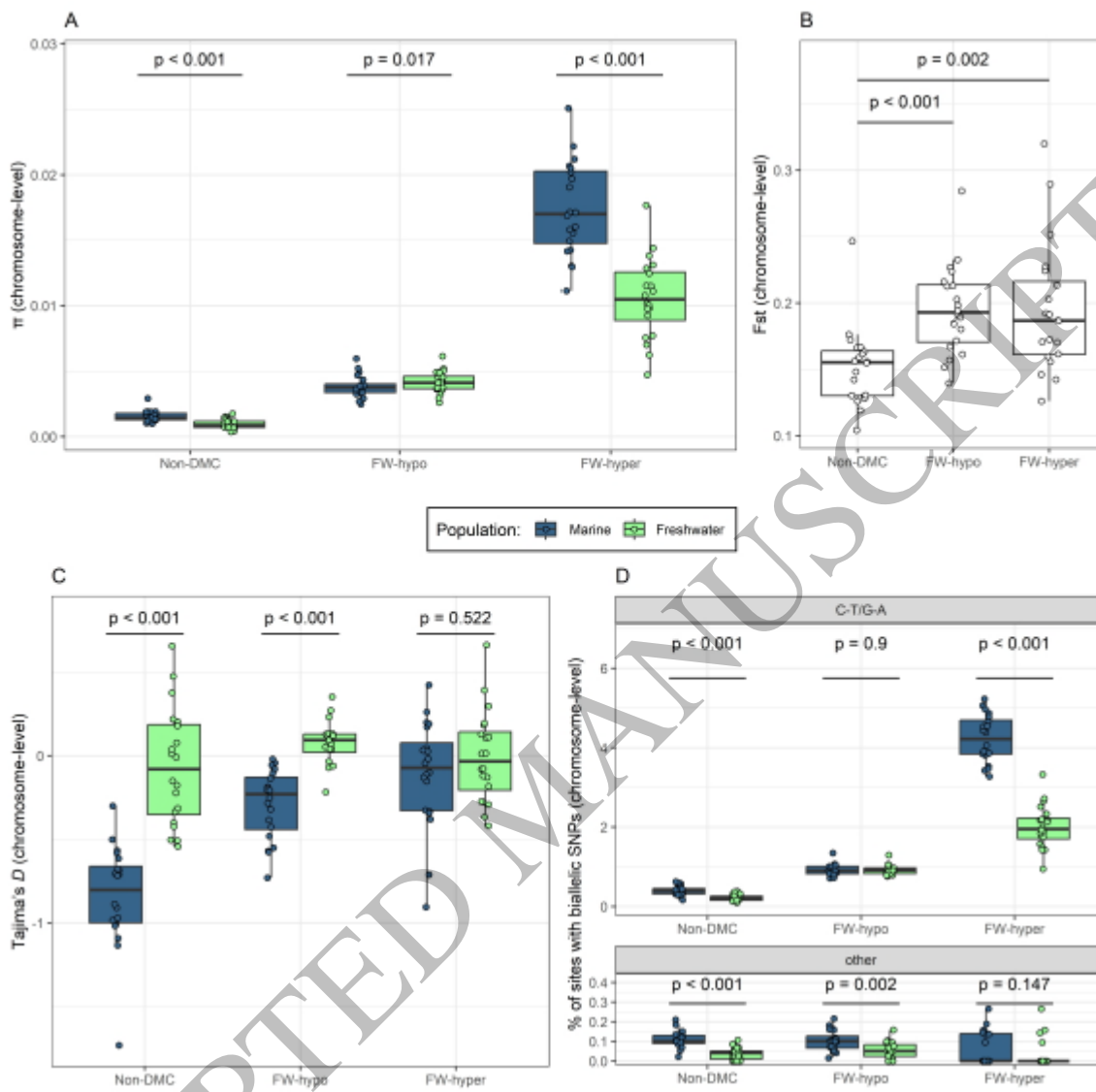


Figure 2
153x150 mm (.97 x DPI)

1
2
3
4

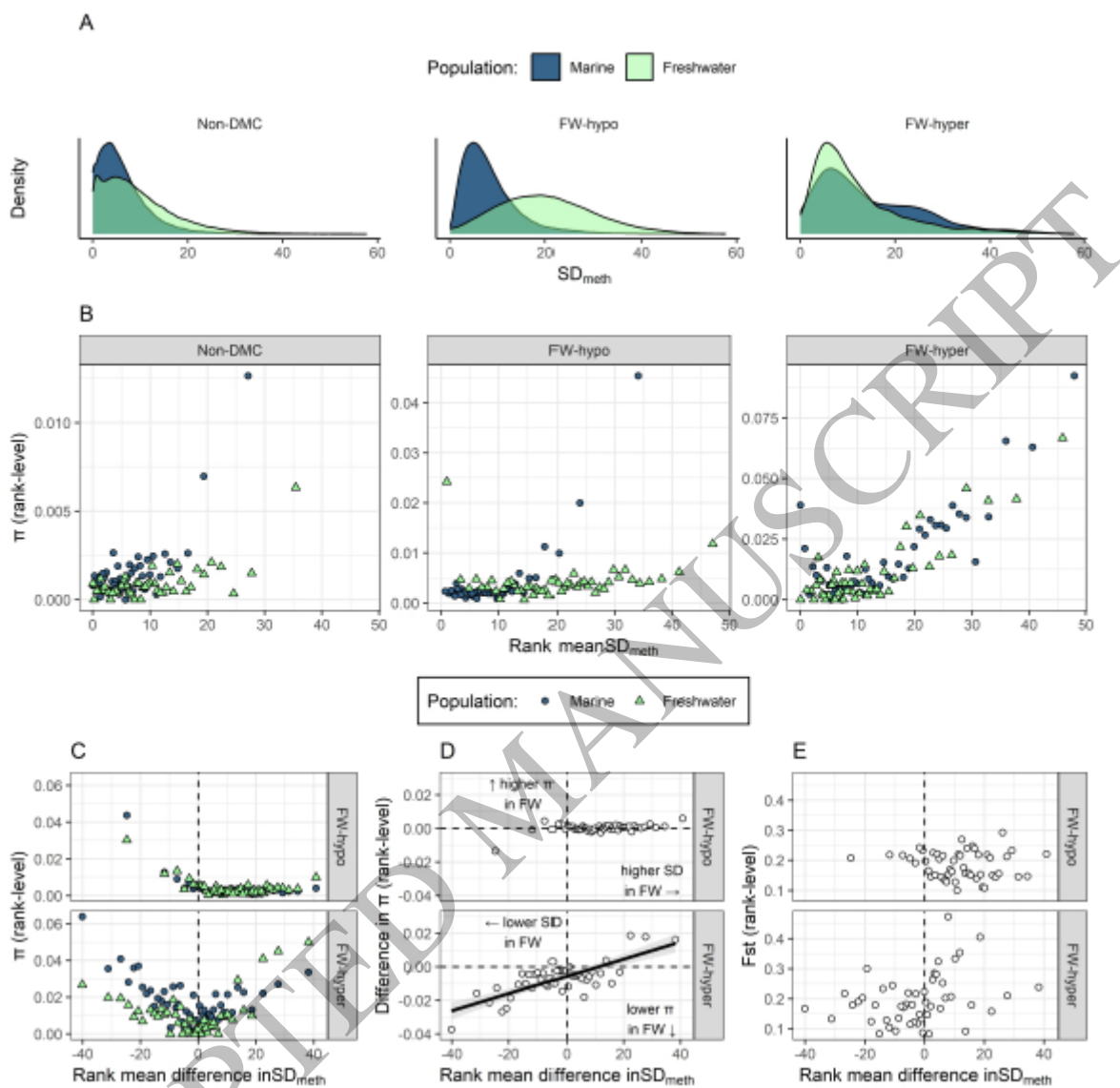


Figure 4
152x152 mm (.97 x DPI)

1
2
3
4

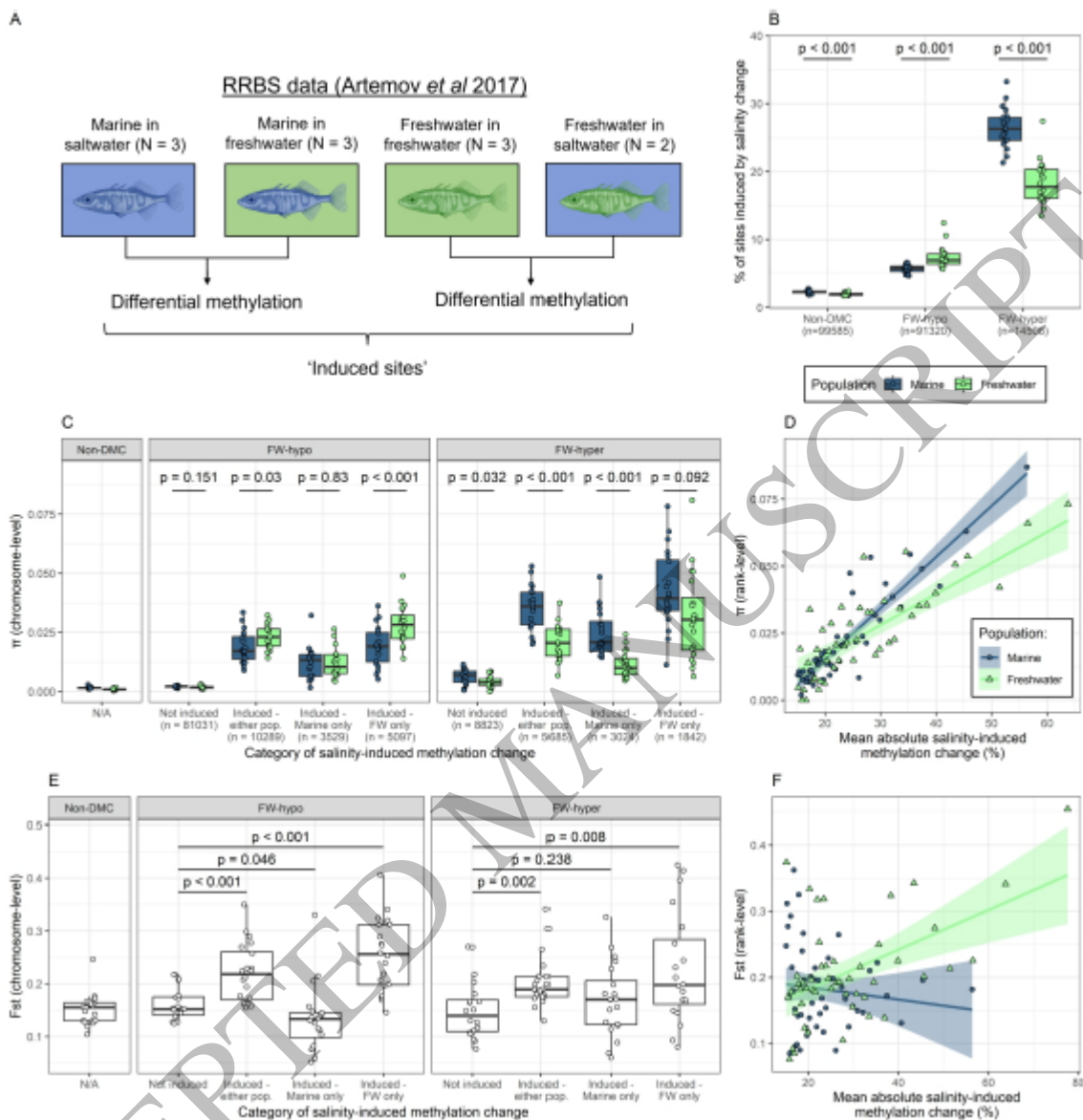


Figure 5
152x156 mm (.97 x DPI)

1
2
3