**Location of SNPs within colony**

**Location of SNPs between genome and GBS reference**

positive values are with 92bp reference, negative are with 90bp reference
gray shows all SNPs, red shows A->C

1

2

3  Supplemental Figure 1. SNP Profiles

4       All of the entries in the reference we extracted from the STACKs pipeline are 92

5  bases long. This histogram shows the number of SNPs per base position within the animals

6  we sequenced (top panel) and when compared with the Lucigen genome (bottom panel). The

7  final two bases show unusually high numbers of SNPs in both cases, and especially high

8  numbers of A to C transversions (red).

9       We found that these excess SNPs are always associated with a restriction cutsite. And

10 specifically, they are associated with regions where the cutsite is within two bases from the

11 end of the reference. That is, all the spurious cases ended with either CGWCX or CGWCXX

12 where 'X' represents the SNP and CGWC is the restriction recognition sequence. Given the

13 way in which GBS libraries are made, we believe that these SNPs may be part of the barcode

14 and/or Illumina adapter sequence that was too short to be recognized and trimmed by BGI.

15       Our first attempt to deal with this issue was by cropping the reads entering the

16    STACKs pipeline to 90 bases as we expected that would remove the short bits of adapter that

17    remained. We tested this by re-running the STACKs pipeline using the new 90bp reference.

18    Unfortunately, the issue was not resolved (see the negative values in both panels above). We

19    found that the easiest way to remove these tailing SNPs was to use the 92bp reference, but

20    blacklist the final two bases from the analyses of diversity (this was implemented in the

21    populations script through a whitelist of SNPs in the first 90 bases rather than a blacklist of

22    SNPs in the final 2 bases because populations does not allow SNP-specific blacklists).

23

24       Despite our solution, the question remained: Why do we find elevated substitution

25    rates near the ends of the reads regardless of the length of the read? We think the answer lies

26    in the fact that alignments only allow so many mismatches. For example, there are reads that

27    have a cutsite four bases from the end of the read (CGWCXXXX), but four mismatched

28    bases is too many to form an alignment (the tailing bases would be soft-clipped rather than

29    aligned as a mismatch) and so these were never analyzed for SNPs (we counted SNPs only in

30    alignments without any clipping). When we later trimmed the tailing two bases off, these

31    cases now had only two mismatches in the alignment, and so passed the alignment-length

32    filter, and the mismatches were identified as SNPs. As cutsites do occur randomly in the

33    genome, it does not matter if we crop two or ten bases from the raw reads, a cutsite would

34    still lie within two bases on some reference somewhere. The best way perhaps to solve the

35    issue would be to do a thorough bout of adapter trimming followed by cropping five to six

36    bases. Unfortunately, we had the libraries built and sequenced by BGI who used proprietary

37    barcodes and adapter sequences. As they were unwilling to share these sequences with us,

38    further cleaning of the reads was not possible and the blacklisting approach that we used

39    seems best.