

# Additional File 1: Calculation of expected SNP frequencies in out- and back-crosses.

Graham J Etherington<sup>1,2†</sup>, Jacqueline Monaghan<sup>1†</sup>, Cyril Zipfel<sup>1</sup> and Dan MacLean<sup>1\*</sup>

\*Correspondence:

dan.maclea@tsl.ac.uk

<sup>1</sup>The Sainsbury Laboratory,  
Norwich Research Park, NR4  
7UH, Norwich, UK

Full list of author information is  
available at the end of the article

<sup>†</sup>Equal contributor

## SNP frequency in different crosses

Methods such as SHOREMap [1] and NGM [2] rely on the relative density of SNPs of one type to SNPs of another type to work (usually the number of homozygous SNPs relative to heterozygous SNPs). As a second ecotype is involved in an out-cross then you get many thousands of SNPs introduced. We can define some terms

- $S_{\text{ler}}$  the number of SNPs between the reference sequence and another ecotype (e.g Ler)
- $S_{\text{ref}}$  the number of SNPs between the reference sequence (i.e Col-0) and the bulked mutants (i.e *bak1-5 mob1* and *bak1-5 mob2*).
- $S_{\text{parent}}$  the number of SNPs present in the parent (i.e in *bak1-5*)
- $S_{\text{mut}}$  the number of SNPs present in the mutant

The formula for the number of SNPs in an out-cross is

$$S_{\text{ref}} \cup S_{\text{parent}} \cup S_{\text{ler}} \quad (1)$$

(where the symbol  $\cup$  is the symbol that defines the 'union' operation, the joining of two sets with removal of overlaps). So this formula represents the non-redundant, summed set of the SNPs in each component of the background. Looking at the relative magnitudes of the numbers in each of those subsets of SNPs

- $S_{\text{ref}}$  is on the order of 1200 [3]
- $S_{\text{parent}}$  is on the order of 500 [4]
- $S_{\text{mut}}$  is on the order of 500 for the same reasons as above
- $S_{\text{ler}}$  is on the order of 150,000 [1, 3]

So in the Arabidopsis genome we would expect 1 SNP every 900 nt in a out-crossed experiment with another ecotype and 1 SNP every 65,000 nt in the back-cross.

We see then that any cross with another ecotype is going to be more amenable to statistical methods. The number of SNPs from a backcross is not amenable to frequentist approaches over an entire genome.

## Author details

<sup>1</sup>The Sainsbury Laboratory, Norwich Research Park, NR4 7UH, Norwich, UK. <sup>2</sup>The Genome Analysis Centre Norwich Research Park, NR4 7UH, Norwich, UK.

## References

1. Schneeberger, K., Ossowski, S., Lanz, C., Juul, T., Petersen, A.H., Nielsen, K.L., Jørgensen, J.-E., Weigel, D., Andersen, S.U.: SHOREmap: simultaneous mapping and mutation identification by deep sequencing. *Nature methods* **6**(8), 550–551 (2009)
2. Austin, R.S., Vidaurre, D., Stamatidou, G., Breit, R., Provart, N.J., Bonetta, D., Zhang, J., Fung, P., Gong, Y., Wang, P.W., McCourt, P., Guttman, D.S.: Next-generation mapping of Arabidopsis genes. *The Plant journal : for cell and molecular biology* **67**(4), 715–725 (2011)
3. Ossowski, S., Schneeberger, K., Clark, R.M., Lanz, C., Warthmann, N., Weigel, D.: Sequencing of natural strains of Arabidopsis thaliana with short reads. *Genome research* **18**(12), 2024–2033 (2008)
4. Weigel, D., Glazebrook, J.: EMS Mutagenesis of Arabidopsis Seed. *Audio and Electroacoustics Newsletter, IEEE* **2006**(5), (2006)