Supplemental Materials

Removing reference mapping biases using limited or no genotype data identifies allelic differences in protein binding at disease-associated loci

Martin L. Buchkovich, Karl Eklund, Qing Duan, Yun Li, Karen L. Mohlke and Terrence S. Furey

Supplemental Figure Legends

Figure S1. Reference mapping biases influence sequence alignment. When mapping reads to a reference genome containing a single allele at heterozygous sites (left), sequence reads containing the reference allele (C, blue) are more likely to map correctly than reads containing the non-reference allele (A, red), especially in the presence of sequencing errors (orange). Reads containing each allele are equally as likely to map correctly when mapping reads to sequence containing both alleles at heterozygous sites (right).

Figure S2. False positive imbalance sites have more significant p-values. Boxplot of *P*-values at sites of allelic imbalance using complete (left), and partial (middle) genotypes and common variants (right). For the complete genotype alignment, *P*-values are further subdivided by inclusion in partial genotypes and common variants. For partial genotypes and common variant alignments, all *P*-values are displayed in addition to being divided into inclusion during alignment and predicted from the alignment, and then further divided into whether or not the sites was predicted using complete genotypes (false positive vs true positive). Subdivided groups with significant differences in *P*-value distributions (Mann Whitney U test; *P*<.01 and *P*<.001) are indicated.

Figure S3. Read length influences sequence alignment and allelic imbalance detection at heterozygous sites. (A) Alignment statistics for alignments of CREB1 ChIP-seq data, and when using different (B) read lengths and (C) sequence depths, plotted as percent of the 50 bp statistics. (D) Histogram of the number of reads containing the underrepresented allele at sites with significant imbalance (binomial P<.01, uncorrected) with 2 or more reads containing each allele. Vertical dashed line indicates the minimum of 5 or more reads containing each allele required for imbalance detection in (A).

Figure S4. Correlation of alignment statistics and number of imbalances detected. The number of sites of allelic imbalance in thirteen ChIP-seq and one DNase-seq dataset compared to the (A) total number of reads aligned; (B) number of reads aligned to heterozygous sites; (C) total number of bases aligned (read length x total reads aligned); the percent of genome with greater than (D) 1X and (E) 10X coverage; (F) the average read depth at bases with 1X or more coverage; (G) the ratio of sites with 10X coverage to 1X coverage and (H) the number of heterozygous sites identified. Pearson correlation R² values for each statistic and number of allelic imbalances are displayed for all data and only ChIP-seq data.

Figure S5. Allelic differences in binding at sites without predicted allelic imbalance. EMSA using purified CREB1 and labeled probes containing each allele at five sites with reads

mapping but not significant allelic imbalance to test for allelic differences in binding. Alleles colored blue are the reference allele and alleles colored red are the other allele. Allelic differences in protein binding were detected at two sites (starred) without predicted imbalance but with predicted allelic differences in the presence of CREB1 motif. Only CREB1-bound probe is shown.

Alignment to single allele Alignment to both alleles

Not Aligned 2 or more mismatches	GACCTCTGAAGCAATTA GGCCTCTGAAGCAAGTA	
Aligned 0 or 1 mismatches	GGCGTCTGCAGCAATTA GGCCTCTGCAGCTATTA GGCCTCTGAAGCAATTA GGCCTCTGCAGCAATTA	GACCTCTGAAGCAATTA GGCCTCTGAAGCAAGTA GGCCTCTGCAGCAATTA GGCCTCTGCAGCTATTA GGCCTCTGAAGCAATTA GGCCTCTGCAGCAATTA
Reference Sequence	GGCCTCTGCAGCAATTA f Heterozygous Site C/A	GGCCTCTGCAGCAATTA GGCCTCTGAAGCAATTA T Heterozygous Site C/A





Figure S3



	Motif predicted for			Motif predicted for both alleles					Allele 1		e 2	Imbalance	
	reference allele only		Variant			rsID	Allele	Reads	Allele	Reads	P-value		
	Δ	*	B*	С	D	F	A	rs72807213	G	24	Α	17	0.35
	G	A	TC	GA	TA	GA	В	rs9388486	Т	36	С	33	0.81
Purified CREB1	-		-				С	rs274035	G	19	Α	17	0.87
	•					D	rs12145434	Т	78	Α	82	0.81	
	-						E	rs55811458	G	42	A	53	0.30