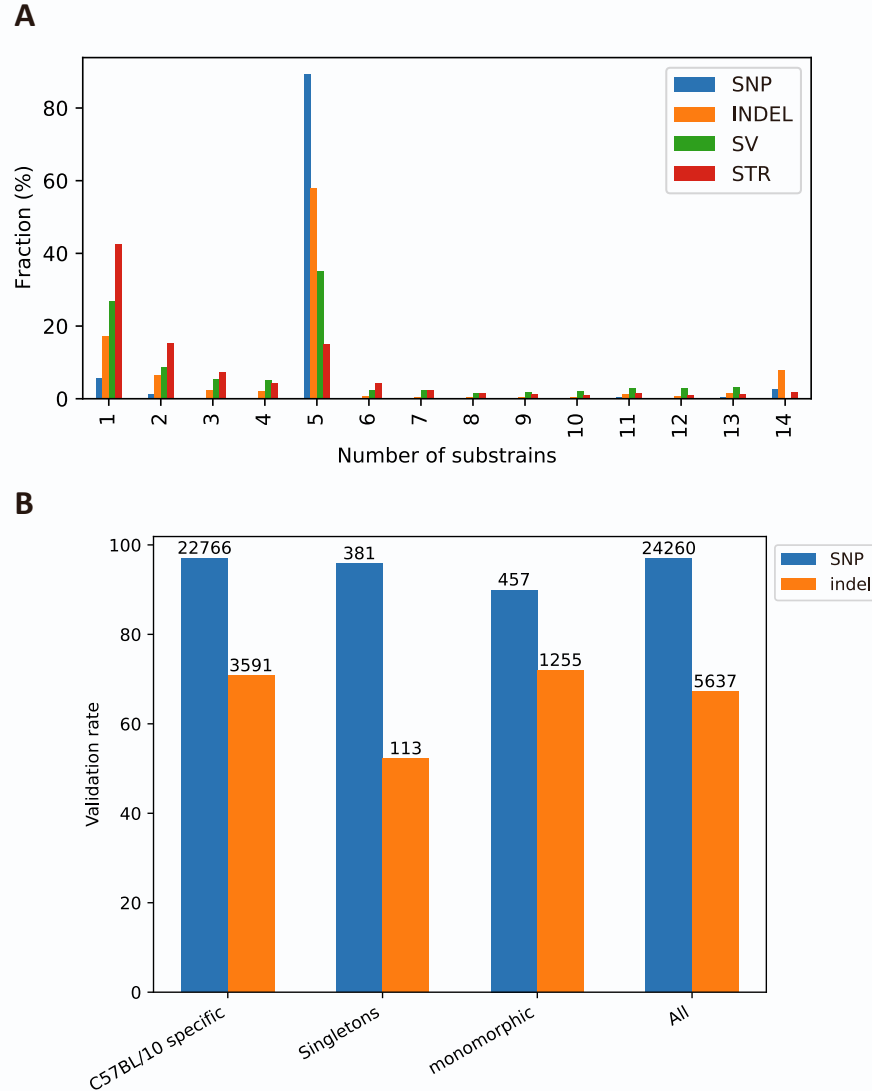


**Cell Genomics, Volume 2**

**Supplemental information**

**SNPs, short tandem repeats, and structural variants  
are responsible for differential gene expression  
across C57BL/6 and C57BL/10 substrains**

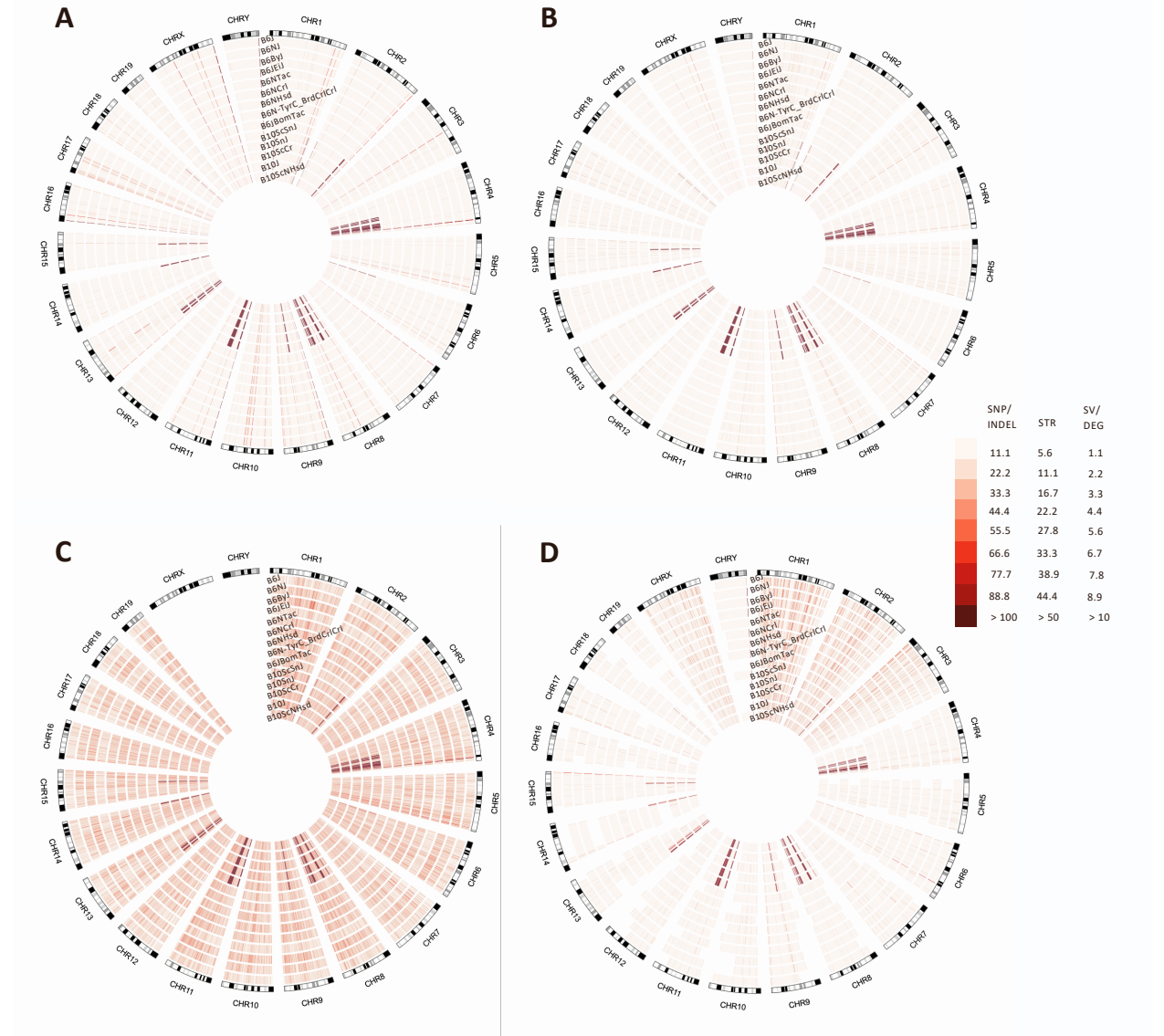
**Milad Mortazavi, Yangsu Ren, Shubham Saini, Danny Antaki, Celine L. St. Pierre, April Williams, Abhishek Sohni, Miles F. Wilkinson, Melissa Gymrek, Jonathan Sebat, and Abraham A. Palmer**

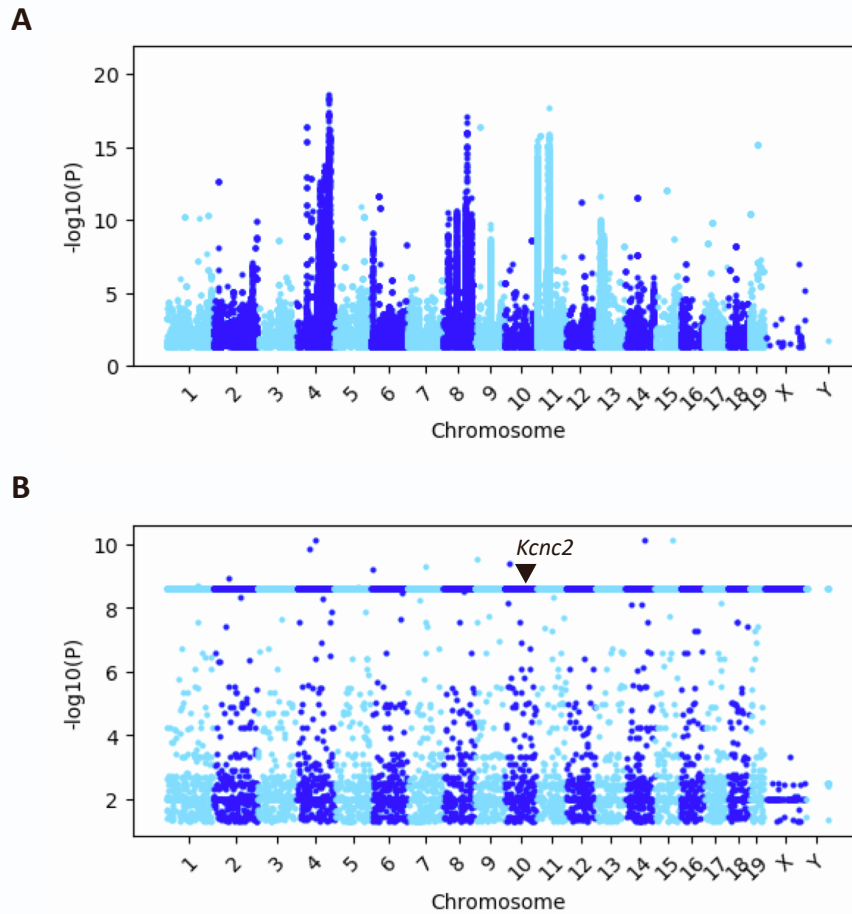


**Figure S1. Variant frequency and validation rates, related to STAR Methods. A:** Fraction of variants observed in substrains (five C57BL/10 and nine C57BL/6 substrains) in each variant category. The spike at 5 reflects polymorphisms that separated C57BL/10 (n=5) from C57BL/6 (n=9) substrains. The smaller spike at 14 represents instances where none of the substrains (including C57BL/6J, which is the basis for mm10) matched the mm10 reference genome. **B:** Validation rates of WGS variants in the protein coding regions using RNA-Seq data. WGS SNPs and INDELS which intersect with protein coding exon and UTR annotations (from Ensembl) and have at least 3X coverage in RNA-Seq dataset are considered for validation. Variants from RNA-Seq data were called by GATK best practices[S1] for each substrain separately (see also STAR Methods). Validation rate between different categories of variants are compared. The total number of variants in each category is indicated on top of each bar. Overall 97% of all SNPs and 67% of all INDELS were validated using RNA-Seq data.



segments with C57BL/10-specific SNP hotspots. X-axis shows strains that have at least 300 common loci and at least 90% concordance with C57BL/10-specific SNPs in each segment. The SNP data for the strains is obtained from MGI[S2]. The segments are color coded with the concordance value. The strain labels on the x-axis are color coded with blue: *domesticus* origin, and red: *musculus* origin[S3].





**Figure S4. Association of all genomic variants and expression of DEGenes, related to Figure 2.** Association tests of DEGene expressions of C57BL/6 and C57BL/10 substrains with all genomic variants (SNPs, INDELs, STRs and SVs) was performed by linear regression model with Limix[S4] **A:** Association of DEGene expressions with all variants (SNPs, INDELs, STRs and SVs) in the *cis*-region defined as 1Mb upstream and 1Mb downstream of the DEGene. The p-values are plotted at the genomic locations of the corresponding DEGenes. **B:** Association of *Kcnc2* expression with all genomic variants across the genome shows that variants with the same strain distribution pattern have identical p-values. The flat horizontal line at about  $-\log_{10}(p) = 8.4$  reflects features that have the same strain distribution pattern and therefore all yield identical p-values when tested for association with the gene expression data.

Strain	Substrain	CNVnator		Lumpy		
		DEL	DUP	DEL	DUP	INV
C57BL/6	C57BL/6J	27	461	11	97	2
	C57BL/6NJ	5	456	65	105	2
	C57BL/6ByJ	11	426	82	90	5
	C57BL/6JeiJ	14	469	42	91	4
	C57BL/6Ntac	27	465	64	77	4
	C57BL/6NCrI	9	420	66	92	2
	C57BL/6NHsd	127	488	73	86	3
	B6N-TyrC/BrdCrCrI	7	445	63	96	3
	C57BL/6JbomTac	13	447	50	78	3
C57BL/10	C57BL/10ScSnJ	95	430	1204	129	15
	C57BL/10SnJ	89	420	1195	125	13
	C57BL/10ScCr	269	481	1196	103	15
	C57BL/10J	74	437	1193	115	16
	C57BL/10ScNHsd	72	426	1192	123	14
	TOTAL	448	1308	1369	279	21

**Table S1. Number of SVs found in C57BL substrains, related to Methods details.**







## References:

[S1] Van der Auwera, G.A., and O'Connor, B.D. (2020). *Genomics in the Cloud* (O'Reilly Media, Inc.).

[S2] MGI Search Mouse SNPs. <http://www.informatics.jax.org/snp>.

[S3] Yang, H., Wang, J.R., Didion, J.P., Buus, R.J., Bell, T.A., Welsh, C.E., Bonhomme, F., Yu, A.H., Nachman, M.W., Pialek, J., et al. (2011). Subspecific origin and haplotype diversity in the laboratory mouse. *Nature genetics* 43, 648-655.

[S4] Lippert, C., Casale, F.P., Rakitsch, B., and Stegle, O. (2014). LIMIX: genetic analysis of multiple traits. *bioRxiv*, 10.1101/003905.